



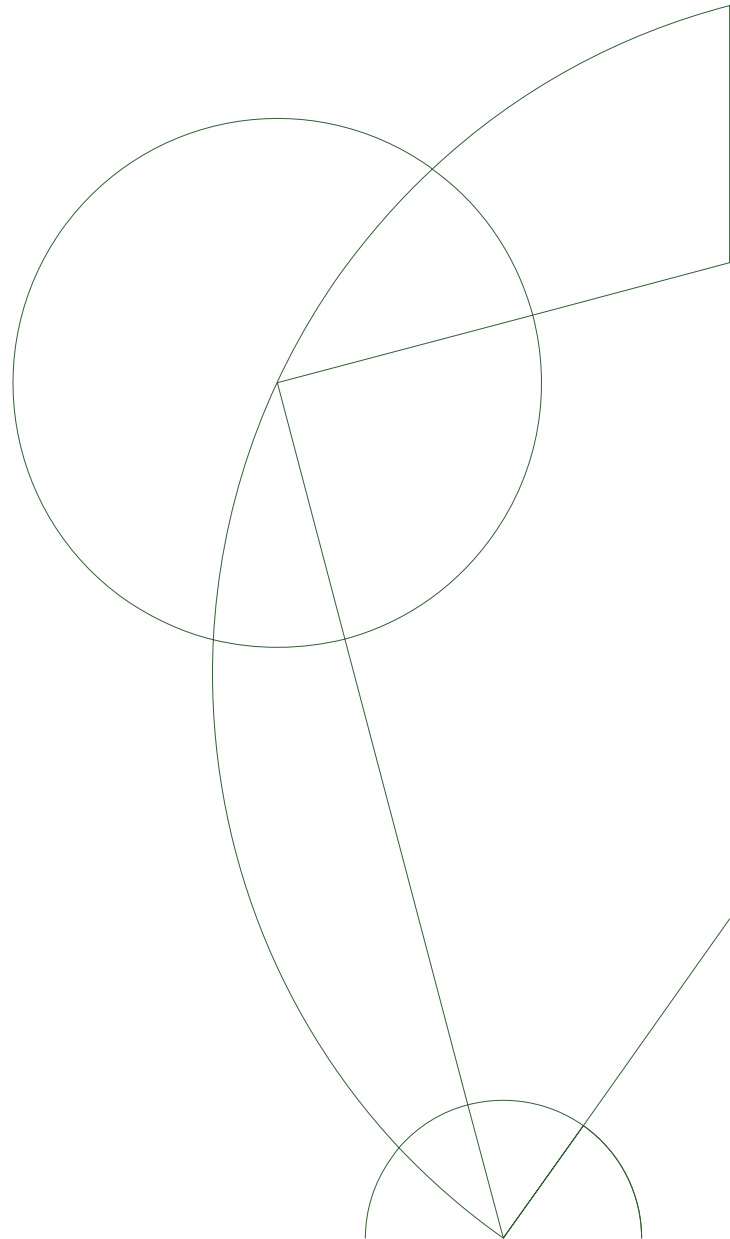
# Thesis of MSc in Mathematics

## Unsupervised bilingual dictionary induction with stable GANs

**Author:** Xuwen Zhang <dlv618@alumni.ku.dk>

**Supervisor:** Anders Søgaard <soegaard@di.ku.dk>

**Submitted on:**



---

## Abstract

# Contents

<b>Contents</b>	<b>3</b>
<b>1 Introduction</b>	<b>4</b>
1.1 Background . . . . .	4
1.2 Vanilla GANs and it's improved versions . . . . .	5
<b>2 Vanilla GAN</b>	<b>8</b>
2.1 working mechanism . . . . .	8
2.2 Issues . . . . .	11
<b>3 Wasserstein GAN</b>	<b>17</b>
3.1 Several probability metrics . . . . .	17
3.2 Why Wasserstain GAN is better . . . . .	19
3.3 Wasserstein GAN . . . . .	23
<b>4 Several improved GANs</b>	<b>25</b>
4.1 Wasserstein GAN with Gradient Penalty (WGAN-GP) . . . . .	25
4.2 Wasserstein GAN with A Consistency Term (CT-GAN) . . . . .	27
<b>5 Experiments</b>	<b>29</b>
5.1 data description . . . . .	29
5.2 Multilingual Unsupervised and Supervised Embeddings . . . . .	29
5.3 Experimental set up . . . . .	32
5.4 Results . . . . .	34
<b>6 Conclusions</b>	<b>40</b>
6.1 Future work . . . . .	40
<b>A Some deep learning techniques</b>	<b>41</b>
A.1 Layer normalization . . . . .	41
A.2 RMS prop . . . . .	42
A.3 Adam . . . . .	42
A.4 leaky Relu . . . . .	43
<b>Bibliography</b>	<b>44</b>

# Chapter 1

## Introduction

### 1.1 Background

Word embedding is a representation for vocabulary in a document, more specifically, it represents pieces of vocabulary as vectors. Word embedding is able to capture information about how different words are related to each other based on the context in a document. The most popular and powerful technique to learn word embeddings is called Word2Vec which is a shallow neural network, and it was introduced by [Mikolov et al. \(2013\)](#).

Cross-lingual word embeddings means mapping a source monolingual word embedding to another target embedding space in order to align both languages together in the same shared vector space. By [Ruder et al. \(2017\)](#), Cross-lingual word embeddings attract NLP researchers for two reasons:

- It enables computation of cross-lingual word similarities.
- Knowledge can be transferred between different languages, even transfer resource-rich language to resource-insufficient language like English to Finnish.

Broadly speaking, there are two methods to learn cross-lingual embeddings, supervised method and unsupervised method. [Conneau et al.\(2017\)](#) reviewed the supervised method in the beginning of their paper, then contributed a new unsupervised way to learn cross-lingual embeddings. The supervised method assumed that for a given pair of language embeddings, a small dictionary is given such that it would be used as an anchor point to align the two language embeddings. More specifically, assuming the embedding dimension of both embeddings is  $d$ ,  $X$  and  $Y$  are matrices with size  $d \times n$ , and the known dictionary contains  $n$  pairs of words  $\{x_i, y_i\}_{i \in 1, 2, \dots, n}$  for  $x_i \in X$  and  $y_i \in Y$ , it can be used to learn an optimal linear transformation  $W$  between  $X$  and  $Y$  such that

$$W^* = \underset{W}{\operatorname{argmin}} \|WX - Y\|_F \quad (1.1)$$

where  $W$  is a  $d \times d$  matrix, and  $\|\cdot\|_F$  denotes the Frobenius norm. After the optimal transformation is obtained, for any source word  $s$ , the translation word

$t$  in target embeddings can be found by searching for the nearest neighbor

$$t = \underset{t}{\operatorname{argmax}} \cos(Wx_s, y_t).$$

The result of optimal  $W$  can be improved by orthogonalizing it. The orthogonal improvement is shown by [Xing et al. \(2015\)](#). Thus, the equation (1.1) can be reformed to the Orthogonal Procrustes problem, and a solution would be reached by solving the singular value decomposition(SVD) of  $YX^T$ .

$$W^* = \underset{W}{\operatorname{argmin}} \|WX - Y\|_F = UV^T,$$

Where  $U\Sigma V^T = \operatorname{SVD}(YX^T)$ , the column vectors of  $U$  and  $V$  are orthonormal, and  $\Sigma$  is a diagonal matrix whose all diagonal elements are positive.

The unsupervised method which [Conneau et al.](#) proposed is like this:

- Assume two vector spaces are approximately isomorphic, i.e their exists an invertible linear mapping between them.
- Learn a rough approximation of mapping  $W$  by adversarial model. In particular, by generative adversarial networks (GANs).
- Apply the learned mapping  $W$  to find word pairs for the most frequent words, and then refine  $W$  by Procrusters. The refined  $W$  would be used to translate all words in the embedding space.
- In the end, apply  $W$  to translate words by searching for the nearest neighbor score, and the nearest neighbor algorithm they used is called cross-domain similarity local scaling(CSLS), which invented by [Conneau et al.](#)

## 1.2 Vanilla GANs and it's improved versions

Generative Adversarial Networks (GANs) are deep neural net architectures that consist of two networks, and GAN's developement is considered a revolutionary advancement in deep learning. The first GAN model proposed by Goodfellow ([Goodfellow et al., 2014](#))[3], it is called vanilla GAN. The idea of GANs was inspired from game theory, in which two models are competing against each other, one is called generator, the other is discriminator, both models will become more and more robust during the training process. More specifically,

- The discriminator is a binary classifier denoted as  $D$  which learns to recognize whether the input is coming from a real dataset or a generator's output. Formally,  $D$  calculates the probability of the sample coming from the generator  $G$  or the real input data for each iteration. Usually the label defined as 1 for real data and 0 for synthetic data.
- The generator denoted as  $G$ , takes random noise from the latent space, and outputs synthetic samples. The training purpose for the generator

is to approximate the distribution of real data as close as possible, at the same time the probability for fooling the discriminator would increase too, and this makes higher the probability that the discriminator would treat the synthetic output from the generator as real

To measure the similarity between the distributions of generated samples and real samples, Jensen-Shannon (JS) divergence is used as the measurement instead of Kullback–Leibler (KL) divergence, because JS divergence is symmetric and smoother than KL divergence.

GANs are able to generate data from scratch based on the given sample. The main application area of GANs is computer vision, e.g., creating animated characters (Jin et al. 2017) which could be applied by game developing or animation production companies, high resolution images production based on given low resolution images (Ledig et al. 2016), and new video sequence generation (Vondrick et al. 2016).

GANs can also be applied to other domains such as music generation (Fedus et al. 2018), text generation (Mogren. 2016), or even text to image synthesis (Reed et al. 2016) etc. These applications show that GANs can be adapted immediately for commercial purpose.

Although vanilla GANs achieved great success, there are several disadvantages which makes vanilla GANs' training procedure unstable,

- Mode collapse happens frequently, which means the generator only learns limited distributions of real data.
- Lack of metric that can show the training process, therefore the training is fully manual.
- Non- convergence, i.e. it is difficult for the model to reach the Nash equilibrium due to oscillating change of gradients.
- Gradient vanished, which means the generator's weights stop updating due to the discriminator performing too well.
- Highly sensitive for tuning hyperparameters.

Several improved versions of GANs attempt to solve the above issues, e.g., Wasserstein GAN (Arjovsky et al. 2017), in which the Earth Mover (Wasserstein) Distance is applied to substitute minimize the JS divergence as the new loss function to make gradient more stable due to the EM distance is much smoother. Although the EM distance has nice properties, it is almost impossible to calculate, therefore, the Kantorovich-Rubinstein duality (Villani. 2008) algorithm is applied, where all weights in the discriminator are constrained to a 1-Lipschitz function, i.e., all weights in the discriminator must be clipped in a bounded interval. However, weight clipping will bring new issues, even Arjovsky pointed out in the original Wasserstein GAN paper that it is a horrible method to enforce a Lipschitz constraint. Instead of weight clipping, the gradient penalty (Gulrajani, et al. 2017) enforces Lipschitz by penalizing the norm of the discriminator's gradient. Our best models are obtained from this

## 1.2. Vanilla GANs and it's improved versions

---

algorithm. Another algorithm is CT-GAN ([Wei, et al. 2018](#)), where a consistency term is added to the WGAN with gradient penalty.

We focused on the unsupervised part of the MUSE in this thesis and do the following:

1. Try to reproduce the results of [Conneau et al](#) for several language pairs, and use the reproduced results as benchmark.
2. Implement the improved GANs which are mentioned previously, and compare our result with the benchmark. However, since we found that the WGAN-GP applied MUSE has the best performance, we focused on it mostly.
3. Analysing the results.

## Chapter 2

# Vanilla GAN

This chapter shows how [Goodfellow et al](#) interpret the mechanism of vanilla GAN, and some of the issues [Arjovsky et al](#) explored.

### 2.1 working mechanism

The training purpose for the generator  $G$  is to make the distribution of its output more and more close to real data, the goal of the discriminator  $D$  is to achieve a Nash equilibrium for which it is impossible to distinguish whether its input comes from real or generated dataset. Formally speaking, the task of the discriminator  $D$  is to maximize

$$\mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log D(\mathbf{x})]$$

which means increasing the probability to identify samples coming from real data, as well as to reduce the probability  $D(G(\mathbf{x}))$  to zero by maximizing

$$\mathbb{E}_{\mathbf{x} \sim p_g(\mathbf{x})} [\log(1 - D(G(\mathbf{x})))].$$

Thus, the objective function of the discriminator is

$$\max_D L(D) = \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_g(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$

At the same time, the task for training the generator  $G$  is to increase the probability of the output of  $D$  to assign labels to fake data, therefore the objective function of the generator is

$$\min_G L(G) = \mathbb{E}_{\mathbf{x} \sim p_g(\mathbf{x})} [\log(1 - D(G(\mathbf{x})))].$$

Since  $D$  and  $G$  are playing a minimax game, the final objective function is

$$\min_G \max_D L(G, D) = \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_g(\mathbf{z})} [(1 - \log D(G(\mathbf{z})))] \quad (2.1)$$

The notation  $p_r(\mathbf{x})$  means the probability distribution for real data,  $p_z(\mathbf{x})$  denotes generated fake data.



**Proposition 2.1:** [ [Proposition 1 \(Goodfellow et al., 2014\)](#) ]

For fixed  $G$ , the optimal value for  $D$  is:

$$D^*(\mathbf{x}) = \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})}. \quad (2.2)$$

*Proof:*

Here is a summary of the proof by the original author:

Since the loss function in continuous probability distribution is

$$L(G, D) = \int_{\mathcal{X}} \left( p_r(\mathbf{x}) \log(D(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D(\mathbf{x})) \right) d\mathbf{x}$$

the maximum value of  $L(D, G)$  can be achieved by solving the function inside the integral. For convenience, define the function inside the integral as

$$f(y) = p_r(x) \log y + p_g(x) \log(1 - y),$$

and take the derivative of loss with respect to  $x$  and set it to zero, then get

$$\begin{aligned} \frac{df(y)}{dy} &= p_r(x) \frac{1}{y} - p_g(x) \frac{1}{1-y} \\ &= \frac{p_r(x) - y(p_r(x) + p_g(x))}{y(1-y)} \\ &= 0 \\ \implies y &= \frac{p_r(x)}{p_r(x) + p_g(x)} \end{aligned}$$

Apply the above to a vector variable  $x$ .

Thus, the  $D^*(\mathbf{x}) = \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \in [0, 1]$ . □

By proposition 1 above, The equation (2.1) can be rewritten as

$$\begin{aligned} L(G, D^*) &= \max_D L(G, D) \\ &= \mathbb{E}_{\mathbf{x} \sim p_r} [\log D^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [(1 - \log D^*(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_r} \left[ \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[ \frac{p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] \end{aligned}$$

**Theorem 2.1:** [ [Theorem 1 \(Goodfellow et al., 2014\)](#) ]

The global minimum of  $L(G, D^*)$  can be reached if and only if  $p_r(\mathbf{x}) = p_g(\mathbf{x})$ , at this point,  $C(G) = -2\log 2$  and  $D^*(\mathbf{x}) = 0.5$ .

*Proof:*

Here is a summary of the proof by the original author:

By equation (2.2), if  $p_r(\mathbf{x}) = p_g(\mathbf{x})$ , then  $D^*(\mathbf{x}) = 0.5$ , substituting this value to equation (4) we have

$$\begin{aligned}\min_G L(G, D^*) &= \int_{\mathcal{X}} \left( p_r(\mathbf{x}) \log \frac{1}{2} + p_g(\mathbf{x}) \log \frac{1}{2} \right) dx \\ &= \log \frac{1}{2} \left( \int_{\mathcal{X}} p_r(\mathbf{x}) dx + \int_{\mathcal{X}} p_g(\mathbf{x}) dx \right) \\ &= 2 \log \frac{1}{2} \\ &= -2 \log 2.\end{aligned}$$

to show  $-2 \log 2$  is the optimal value for  $L(G, D^*)$ , we need to prove they are equal to each other. Subtracting one from the other and by applying proposition 1 we get

$$\begin{aligned}L(G, D^*) - (-2 \log 2) &= 2 \log 2 + \int_{\mathcal{X}} \left( p_r(\mathbf{x}) \log \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right. \\ &\quad \left. + p_g(\mathbf{x}) \log \frac{p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right) dx \\ &= \left( \log 2 + \int_{\mathcal{X}} p_r(\mathbf{x}) \log \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} dx \right) \\ &\quad + \left( \log 2 + \int_{\mathcal{X}} p_g(\mathbf{x}) \log \frac{p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} dx \right) \\ &= \left( \int_{\mathcal{X}} p_r(\mathbf{x}) \log 2 dx + \int_{\mathcal{X}} p_r(\mathbf{x}) \log \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} dx \right) \\ &\quad + \left( \int_{\mathcal{X}} p_g(\mathbf{x}) \log 2 dx + \int_{\mathcal{X}} p_g(\mathbf{x}) \log \frac{p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} dx \right) \\ &= \int_{\mathcal{X}} p_r(\mathbf{x}) \log \frac{2 \cdot p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} dx + \int_{\mathcal{X}} p_g(\mathbf{x}) \log \frac{2 \cdot p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} dx \\ &= KL\left(p_r \parallel \frac{p_r + p_g}{2}\right) + KL\left(p_g \parallel \frac{p_r + p_g}{2}\right) \\ &= 2JS(p_r \parallel p_g)\end{aligned}$$

Since the Jensen-Shannon divergence equals to zero when the two probability distributions  $p_r$  and  $p_g$  are equal everywhere, then, the global optimal value for the loss is  $-2 \log 2$  and  $D^*(\mathbf{x}) = 0.53227$ .  $\square$

Below is the architecture draft of vanilla GAN and the training algorithm in the original paper (Goodfellow et al., 2014)[3].

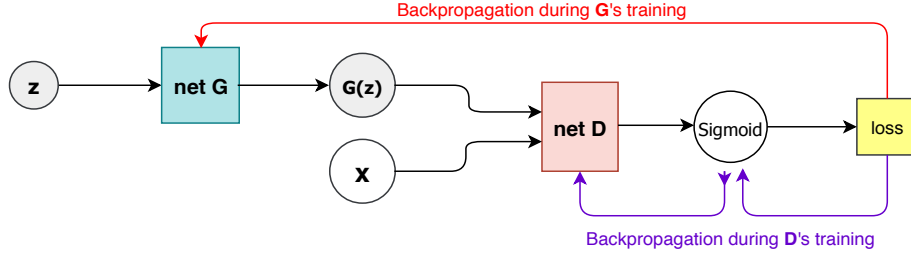


Figure 1: The structure of Vanilla GAN

---

**Algorithm 2.1:** Minbatch training algorithm uses stochastic gradient descent as optimizer,  $n$  is a hyperparameter which for each generator's iteration, trains discriminator  $D$  for  $n$  iterations. The original paper used  $n = 1$ .

---

- 1: **for**  $i = 1, 2, \dots$ , number of training iterations **do**
- 2:   **for** training discriminator  $n$  iterations **do**
- 3:     Sample  $k$  minibatch random noise from latent space distribution  $p_g(z)$
- 4:     Sample  $k$  minibatch input examples from data distribution  $p_r(x)$
- 5:     Backpropagating all weights in  $D$  by maximizing the loss:

$$\nabla_{W_D} \frac{1}{k} \sum_{i=1}^k \left[ \log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right].$$

- 6:   **end for**
- 7:   Sample  $k$  minibatch random noise from latent space distribution  $p_g(z)$
- 8:   Backpropagating all weights in  $G$  by minimizing the loss:

$$\nabla_{W_G} \frac{1}{k} \sum_{i=1}^k \log (1 - D(G(z^{(i)}))).$$

- 9: **end for**
- 

However, [Goodfellow et al](#) pointed out that in practice, when the training has just started, the distributions of  $\mathbb{P}_r$  and  $\mathbb{P}_g$  may have a huge difference, in this case,  $D$  can easily distinguish whether the data comes from training data or synthetic. Therefore,  $\log(1 - D(G(z)))$  saturates. Maximizing  $\log D(G(z))$  can provide more robust gradients in the early training instead of minimizing  $\log(1 - D(G(z)))$ .

## 2.2 Issues

By theorem 1, the minimum of the discriminator's loss is

$$L(G, D^*) = 2JS(\mathbb{P}_r || \mathbb{P}_g) - 2\log 2$$

However, the discriminator's loss will converge to zero in practice . Numerically,  $2\log 2 = 1.386294$ , and this means that  $JS(\mathbb{P}_r || \mathbb{P}_g)$  will converge to 0.693 instead of zero, and at this point the distributions of  $\mathbb{P}_r$  and  $\mathbb{P}_g$  are very similar.

By [Arjovsky et al \[20\]](#), the reason which cause the situation above happens is that either both distributions have disjoint supports, or the distributions are discontinuous, in other words, probability density function does not exist. [Arjovsky et al \[20\]](#) clarified that the term continuous means absolutely continuous. More specifically, for an absolutely random variable, it has the property that if  $P(X \in B) = 0$  then  $B$  has 0 Lebesgue measure. Whereas a random variable is continuous if  $P(X = x) = 0 \quad \forall x \in X$ , where  $x$  is a point. Absolutely continuous implies continuous since the Lebesgue measure of points is 0, but not vice versa. Moreover, If a random variable has support on low dimensional manifold, it is not absolutely continuous. [Narayanan et al. \(2010\)](#) claimed that the support of distribution of  $\mathbb{P}_r$  lies on low dimensional manifold.

For  $\mathbb{P}_g$ , if  $\dim(\mathcal{Z}) < \dim(\mathcal{X})$ , then  $\mathbb{P}_g$  is not continuous. This concept is mathematically formalized in the following lemma.

**Lemma 1:** [ [Lemma 1. \(Arjovsky et al. 2017\)](#) ] Assume  $g : \mathcal{Z} \mapsto \mathcal{X}$  be a neural network, which is a function composed by affine transformations with non-linear activation functions, then a countable union of manifolds which contains  $g(\mathcal{Z})$  would have dimensionality less than or equal to  $\dim(\mathcal{Z})$ . Thus, the measure of  $g(\mathcal{Z})$  would be zero in  $\mathcal{X}$  if  $\dim(\mathcal{Z}) < \dim(\mathcal{X})$ .

Since generator  $G$  is a neural network, this lemma states that if the input dimension of  $z$  less than  $g(z)$ , then the learning ability of  $g(z)$  will be limited, only small part of input distribution can be learned by  $g(z)$  due to the zero measure of  $g(z)$ .

**Lemma 2:** (Smooth Urysohn's Lemma ) [ [Urysohn's lemma](#) ]

Let  $\mathcal{M}$  and  $\mathcal{S}$  be two compact and disjoint subsets in a normal topological space  $(\mathcal{X}, \mathcal{T})$ , then there exists a function  $f : \mathcal{X} \mapsto [0, 1]$  which is smooth such that  $f(x) = 0 \quad \forall x \in \mathcal{S}$  and  $f(x) = 1 \quad \forall x \in \mathcal{M}$ .

The next theorem state that if the supports of two distributions  $\mathbb{P}_r$  and  $\mathbb{P}_g$  are disjoint, there exists a perfect discriminator can be reached, along with vanishing gradients. In this case, therefore, when  $D$  trained too well, the weights update of  $G$  will finish, and make  $G$  stop to approximate.

**Theorem 3:** [ [Theorem 2.1 \(Arjovsky et al. 2017\)](#) ]

Assume the supports of  $\mathbb{P}_r$  and  $\mathbb{P}_g$  included on disjoint compact subsets  $\mathcal{M}$  and  $\mathcal{S}$  respectively, then there exists a complete accurate smooth discriminator  $D^* : \mathcal{X} \mapsto [0, 1]$  which have 0 error, plus  $\nabla_x D^*(x) = 0, \quad \forall x \in \mathcal{M} \cup \mathcal{S}$ .

*Proof:*

By assumption,  $\mathcal{M}$  and  $\mathcal{S}$  are disjoint and compact. Supposed  $d(\mathcal{M}, \mathcal{S}) =$

$\epsilon > 0$  is the distance between  $\mathcal{M}$  and  $\mathcal{S}$ , then define another two sets

$$\begin{aligned}\mathcal{M}' &= \{x : d(x, \mathcal{M}) \leq \frac{\epsilon}{3}\} \\ \mathcal{S}' &= \{x : d(x, \mathcal{S}) \leq \frac{\epsilon}{3}\}\end{aligned}$$

therefore  $\mathcal{M}'$  and  $\mathcal{S}'$  are also both disjoint and compact. Thus, by applying Smooth Urysohn's Lemma, exists a continuous function  $D^* : \mathcal{X} \mapsto [0, 1]$  such that  $D^*|_{\mathcal{S}'} \equiv 0$  and  $D^*|_{\mathcal{M}'} \equiv 1$ . Since  $\forall x \in \text{supp}(D^*|_{\mathbb{P}_r}) = \{x \in \mathbb{P}_r | D^* \neq 0\}$ ,  $\log D^*(x) = 0$ , and  $\log(1 - D^*(x)) = 0 \forall x \in \text{supp}(D^*|_{\mathbb{P}_g})$ , therefore discriminator has zero error and reached complete optimal. In addition, for proving  $D^*(x) = 0$ , first define  $x \in \mathcal{M} \cup \mathcal{S}$ , then for  $x \in \mathcal{M}$ , exist an open ball  $\mathcal{B}(x, \frac{\epsilon}{3}) = \mathcal{B}'$  such that  $D^*|_{\mathcal{B}'}$  is constant. Thus  $\nabla_x D^* \equiv 0$ . For  $x \in \mathcal{S}$ , following the same procedure will done the proof.  $\square$

This theorem says, If the intersection of the supports of  $\mathbb{P}_r$  and  $\mathbb{P}_g$  are empty, then there exists a perfect discriminator which is able to separates  $\mathbb{P}_r$  and  $\mathbb{P}_g$ , and  $D$ 's gradient will be vanished on the union of the two supports. This makes  $G$  stop learning when  $D$  trained too well.

**Definition 1: (Transversal intersection)** [ [Definition 2.2 \(Arjovsky et al. 2017\)](#) ] Assume  $\mathcal{M}$  and  $\mathcal{S}$  are two submanifolds of  $\mathcal{F}$ , and  $\mathcal{F} = \mathbb{R}^n$ . If  $\forall p \in \mathcal{M} \cap \mathcal{S}$

$$T_p\mathcal{M} + T_p\mathcal{S} = T_p\mathcal{F}$$

then we say that  $\mathcal{M}$  and  $\mathcal{S}$  intersect transversally.

$T_p\mathcal{M}$  means the tangent space of  $\mathcal{M}$  at point  $p$ . Moreover, if two submanifolds does not have intersection, then they are transversal.

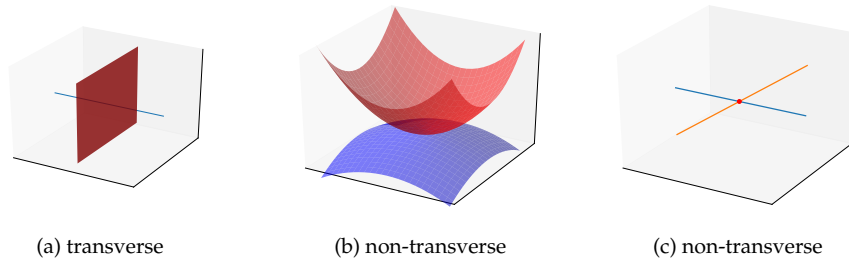


Figure 1: examples of transversal intersection

**Definition 2: (Perfect align)** [ [Definition 2.2 \(Arjovsky et al. 2017\)](#) ] Assume  $\mathcal{M}$  and  $\mathcal{S}$  are two manifolds without boundary. We say  $\mathcal{M}$  and  $\mathcal{S}$  perfectly align if exist  $x \in \mathcal{M} \cap \mathcal{S}$  such that  $\mathcal{M}$  and  $\mathcal{S}$  are non-transversal at  $x$ .

**Lemma 2:** [ [Lemma 2. \(Arjovsky et al. 2017\)](#) ]

Assume  $\mathcal{M}$  and  $\mathcal{S}$  are two regular submanifolds of  $\mathbb{R}^n$ , and they are not full dimension. Supposed  $\epsilon$  and  $\epsilon'$  are two independent continuous random variables, define two perturbed manifolds  $\hat{\mathcal{M}} = \mathcal{M} + \epsilon$  and  $\hat{\mathcal{S}} = \mathcal{S} + \epsilon'$ . Then

$$\mathbb{P}_{\epsilon, \epsilon'}(\hat{\mathcal{M}} \text{ and } \hat{\mathcal{S}} \text{ not perfectly align}) = 1.$$

This lemma means that any small perturbations can make two low dimension manifolds become non-perfect align, i.e., transversally intersect.

**Lemma 3:** [ [Lemma 3. \(Arjovsky et al. 2017\)](#) ]

Assume  $\mathcal{M}$  and  $\mathcal{S}$  are two non-full dimension and non-perfect align regular submanifolds of  $\mathbb{R}^n$ . Let  $\mathcal{M} \cap \mathcal{S} = \mathcal{P}$ .

- If  $\mathcal{M}$  and  $\mathcal{S}$  have boundary, then  $\mathcal{P}$  is a union of no more than four strictly low dimensional manifolds.
- If  $\mathcal{M}$  and  $\mathcal{S}$  without boundary. then  $\mathcal{P}$  is still a manifold which has strictly lower dimension than  $\mathcal{M}$  or  $\mathcal{S}$ .

$\mathcal{P}$  has measure zero in  $\mathcal{M}$  and  $\mathcal{S}$  in both cases.

By lemma 2, two manifolds can be non-perfectly aligned by any subtle disturbance, therefore by lemma 3, their intersection has zero measure.

**Theorem 4:** [ [Theorem 2.2 \(Arjovsky et al. 2017\)](#) ]

Assume  $\mathbb{P}_r$  and  $\mathbb{P}_g$  are two distributions, their supports contained respectively in two closed lower dimensional manifolds  $\mathcal{M}$  and  $\mathcal{S}$  which are non-perfect align. In addition, suppose  $\mathbb{P}_r$  and  $\mathbb{P}_g$  are continuous, which means that for any set  $A \in \mathcal{M}$  with zero measure implies  $\mathbb{P}_r(A) = 0$ , similar for  $\mathbb{P}_g$ . Then there is an full accuracy optimal discriminator  $D^* : \mathcal{X} \mapsto [0, 1]$  which can distinguish  $\mathbb{P}_r$  and  $\mathbb{P}_g$ . Moreover,  $\forall x \in \mathcal{M}$  or  $\forall x \in \mathcal{S}$ ,  $\nabla_x D^*(x) = 0$  and  $D^*$  is smooth in a neighborhood.

In theorem 3, we assumed that the supports of both  $\mathcal{M}$  and  $\mathcal{S}$  are disjoint. However theorem 4 says, even in the condition that the supports of  $\mathcal{M}$  and  $\mathcal{S}$  are join,t and non-perfectly aligned, still exists perfect discriminator which can separate  $\mathcal{M}$  and  $\mathcal{S}$ , which also makes generator's gradient vanished.

**Theorem 5:** [ [Theorem 2.3 \(Arjovsky et al. 2017\)](#) ]

Assume  $\mathbb{P}_r$  and  $\mathbb{P}_g$  are two distributions, their supports contained respectively in two lower dimensional manifolds  $\mathcal{M}$  and  $\mathcal{S}$  which are non-perfect align. Further more, assume  $\mathbb{P}_r$  and  $\mathbb{P}_g$  are continuous in their manifolds. Then

- $JSD(\mathbb{P}_r || \mathbb{P}_g) = \ln 2$
- $KL(\mathbb{P}_r || \mathbb{P}_g) = +\infty$
- $KL(\mathbb{P}_g || \mathbb{P}_r) = +\infty$

This theorem indicates that under the same condition of theorem 4, a perfect discriminator can be learn, then the Jensen Shannon divergence between two

distributions will getting close to  $\ln 2$ . In addition, KL divergence is not available as a metric to evaluate GAN's performance, thus finding a useful metric is the new goal. There are some examples about KL divergence in the next chapter.

**Theorem 6:** [ [Theorem 2.4 \(Arjovsky et al. 2017\)](#) ]

Assume  $g_\theta : Z \mapsto \mathcal{X}$  is a differentiable function which induced the distribution  $\mathbb{P}_g, \mathbb{P}_r$  be the input data distribution,  $D$  is a differentiable discriminator. If the conditions in theorem 2.1 and 2.2 are satisfied,  $\mathbb{E}_{z \sim p(z)} [\|J_\theta g_\theta(z)\|_2^2] \leq M^2$  and  $\|D - D^*\| < \epsilon$ , then

$$\|\nabla_\theta \mathbb{E}_{z \sim p(z)} [\log(1 - D(g_\theta(z)))]\|_2 < M \frac{\epsilon}{1 - \epsilon}$$

*Proof:*

Since the theorem 2.1 and 2.1 satisfied,  $\nabla_x D^*(x) = 0$ . Therefore, by Jensen's inequality and chain rule

$$\begin{aligned} \|\nabla_\theta \mathbb{E}_{z \sim p(z)} [\log(1 - D(g_\theta(z)))]\|_2^2 &\leq \mathbb{E}_{z \sim p(z)} \left[ \frac{\|\nabla_\theta D(g_\theta(z))\|_2^2}{|1 - D(g_\theta(z))|^2} \right] \\ &= \mathbb{E}_{z \sim p(z)} \left[ \frac{\|J_\theta g_\theta(z) \nabla_x D(g_\theta(z))\|_2^2}{|1 - D(g_\theta(z))|^2} \right] \\ &\leq \mathbb{E}_{z \sim p(z)} \left[ \frac{\|\nabla_x D(g_\theta(z))\|_2^2 \|J_\theta g_\theta(z)\|_2^2}{|1 - D(g_\theta(z))|^2} \right] \\ &\leq \mathbb{E}_{z \sim p(z)} \left[ \frac{(\|\nabla_x D^*(g_\theta(z))\|_2 + \epsilon)^2 \|J_\theta g_\theta(z)\|_2^2}{(|1 - D^*(g_\theta(z))| - \epsilon)^2} \right] \\ &= \mathbb{E}_{z \sim p(z)} \left[ \frac{\epsilon^2 \|J_\theta g_\theta(z)\|_2^2}{(1 - \epsilon)^2} \right] \\ &\leq M^2 \frac{\epsilon^2}{(1 - \epsilon)^2} \end{aligned}$$

After square root for both side, get the final result

$$\|\nabla_\theta \mathbb{E}_{z \sim p(z)} [\log(1 - D(g_\theta(z)))]\|_2 < M \frac{\epsilon}{1 - \epsilon}$$

□

**Corollary 1:** [ [Corollary 2.1 \(Arjovsky et al. 2017\)](#) ]

Assume all conditions from theorem 6 are satisfied, then

$$\lim_{\|D - D^*\| \rightarrow 0} \nabla_\theta \mathbb{E}_{z \sim p(z)} [\log(1 - D(g_\theta(z)))] = 0$$

This means generator's gradient is bounded, in other words, the more close the discriminator  $D$  to it's optimal, the small the gradient of generator is. Therefore, in this situation, generator stop learning the input distribution .

In practice, the alternative loss function  $\log D(G(z))$  does not cause gradient vanish, but still it's gradient suffer from unstable update.

**Theorem 7:** [ [Theorem 2.5 \(Arjovsky et al. 2017\)](#) ]

Assume  $\mathbb{P}_r$  and  $\mathbb{P}_g$  are two distributions with density function  $P_r$  and  $P_g$  respectively.  $D^* = \frac{P_r}{P_{g_{\theta_0}} + P_r}$  is the optimal discriminator with a fixed value  $\theta_0$ . Then

$$\mathbb{E}_{z \sim p(z)} \left[ - \nabla_{\theta} \log D^*(g_{\theta}(z)) \Big|_{\theta=\theta_0} \right]$$

The gradient of the new alternative loss function have a inverse KL-divergence and minus JS-divergence between  $\mathbb{P}_r$  and  $\mathbb{P}_g$ . The inverse KL-term can penalize largely different samples by assigns high cost for them, but low cost for mode collapse. Since JS divergence is symmetric, therefore it can not change the problem.



## Chapter 3

# Wasserstein GAN

The way how vanilla GAN works and why it's hard to train already showed up in the previous chapter. [Arjovsky et al. \(2017\)](#) theoretically explored the drawbacks of Vanilla GAN and apply the Earth Mover Distance(EMD) to introduce a new GAN called Wasserstein GAN which tries to solve them, e.g., mode collapse, and lack of interpretable metric that tells us training progress. However, even Wasserstein GAN have better performance than vanilla GAN in image generation, it still have some flaws such weight clipping.

### 3.1 Several probability metrics

There are several formula to measure the similarity between two probability distributions over the same random variable  $x$ . Assumed  $\mathbb{P}_r(\mathbf{x})$  and  $\mathbb{P}_g(\mathbf{x})$  are two distributions.

**Kullback-Leibler (KL) divergence (continuous form):**

$$KL(\mathbb{P}_r||\mathbb{P}_g) = \int_{\mathcal{X}} \log\left(\frac{P_r(\mathbf{x})}{P_g(\mathbf{x})}\right) P_r(\mathbf{x}) dx,$$

it's lower bounded by 0. KL Divergence reaches minimum zero while  $P_r(x) = P_g(x)$  for all  $x \in \mathcal{X}$ , it is asymmetric, . The discrete form of KL divergence is :

$$\begin{aligned} KL(\mathbb{P}_r||\mathbb{P}_g) &= H(p) - H(p, q) \\ &= \sum_x \log\left(\frac{P_r(\mathbf{x})}{P_g(\mathbf{x})}\right) P_r(x) \\ &= - \sum_x \log\left(\frac{P_g(\mathbf{x})}{P_r(\mathbf{x})}\right) P_r(x), \end{aligned}$$

where  $H(p, q)$  is cross entropy between  $p$  and  $q$ . KL divergence is asymmetric.

**Jensen-Shannon(JS) divergence**

$$JS(\mathbb{P}_r, \mathbb{P}_g) = \frac{1}{2} KL(\mathbb{P}_r||\mathbb{P}_m) + \frac{1}{2} KL(\mathbb{P}_g||\mathbb{P}_m)$$

where  $\mathbb{P}_m = \frac{\mathbb{P}_r + \mathbb{P}_g}{2}$ , the JS divergence is symmetric, and bounded in  $[0, 1]$ . It is symmetric, i.e.  $JS(\mathbb{P}_r, \mathbb{P}_g) = JS(\mathbb{P}_g, \mathbb{P}_r)$ .

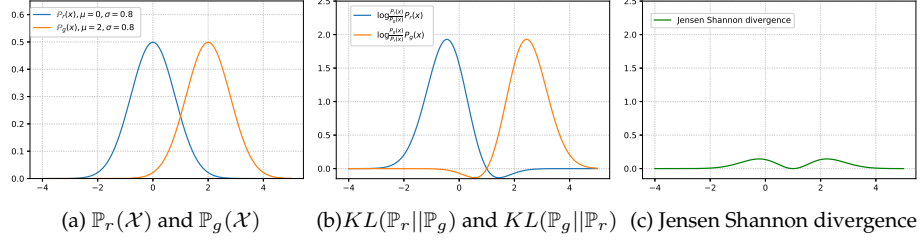


Figure 1: Two different normal distributions with corresponding  $KL$  and  $JS$

#### Earth Mover Distance(EMD):

The last measure is The Earth Mover's Distance (EMD), also called Wasserstein Distance. Informally, assumed there are two different distribution piles of dirt  $\mathbb{P}$  and  $\mathbb{Q}$  with same total amount, the EMD can be interpreted as the minimal cost for moving one distribution piles to another to make both  $\mathbb{P}$  and  $\mathbb{Q}$  have the same distribution. Thus we can calculate EMD as the sum of work flow for moving piles from  $\mathbb{P}$  to  $\mathbb{Q}$ , for each move, the mass times its distance. Sometimes there are lots of ways for move one distribution to another therefore to find the way with minimal cost is an optimization problem.

#### Definition of discrete case:

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \sum_{x, y} \|x - y\| \gamma(x, y) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x, y) \sim \gamma} [\|x - y\|]$$

where  $\Pi(\mathbb{P}_r, \mathbb{P}_g)$  is the set of all combination of joint distributions between  $\mathbb{P}_r$  and  $\mathbb{P}_g$ . Since  $\gamma$  means the percentage of earth should be moved from  $x$  to  $y$ , thus  $\mathbb{P}_g(y) = \sum_x \gamma(x, y)$  and  $\mathbb{P}_r(x) = \sum_y \gamma(x, y)$ , which means  $\mathbb{P}_g(y)$  and  $\mathbb{P}_r(x)$  are marginal distributions of  $\gamma(x, y)$  respectively.

Here is a simple example of how Earth Mover's Distance works, suppose  $X$  and  $Y$  are two different discrete random variables over the spaces  $\mathbb{P}$  and  $\mathbb{Q}$ , then the EMD from  $\mathbb{P}$  to  $\mathbb{Q}$  is

$$\begin{aligned} \text{EMD}(\mathbb{P}, \mathbb{Q}) &= 0.1 * |4 - 3| + 0.1 * |4 - 2| + 0.2 * |4 - 0| \\ &\quad + 0.2 * |5 - 1| + 0.3 * |6 - 2| + 0.1 * |7 - 3| \\ &= 3.5 \end{aligned}$$

### 3.2. Why Wasserstain GAN is better

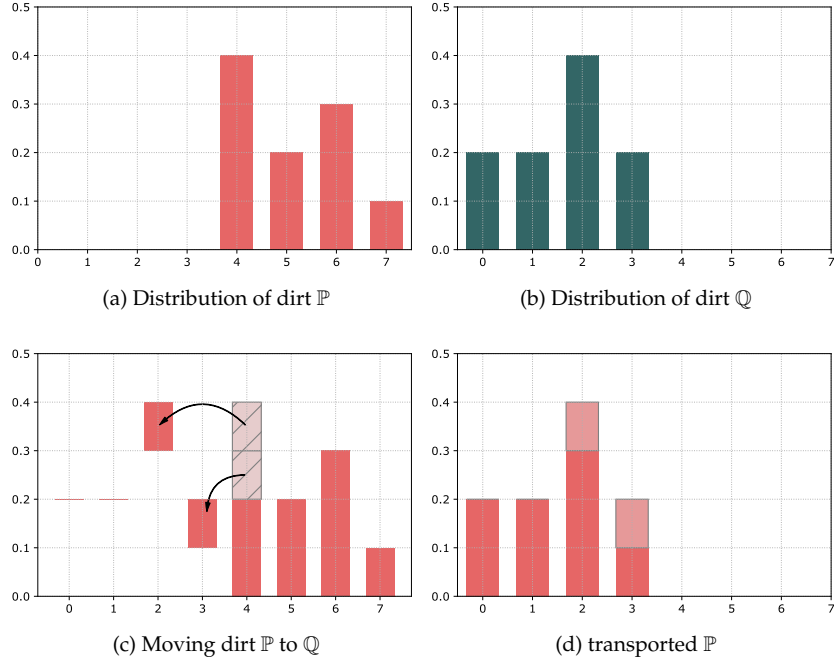


Figure 1: Transporting discrete probability distributions  $\mathbb{P}$  to  $\mathbb{Q}$

### 3.2 Why Wasserstain GAN is better

Here is an example of why using Wasserstein:

Let  $Y$  be the random variable  $Y \sim U[0, 1]$ ,  $\mathbb{P}$  and  $\mathbb{Q}$  are two probability distributions, where  $\mathbb{P}$  is the distribution on  $(0, Y) \in \mathbb{R}^2$ ,  $\mathbb{Q}$  is  $(\theta, Y) \in \mathbb{R}^2$  for  $\theta \in [0, 1]$ . If  $\theta \neq 0$  then we have:

$$\begin{aligned}
 KL(\mathbb{P}||\mathbb{Q}) &= \sum_{x=0, Y \sim U[0,1]} 1 \times \ln \frac{1}{0} = \infty \\
 KL(\mathbb{Q}||\mathbb{P}) &= \sum_{x=\theta, Y \sim U[0,1]} 1 \times \ln \frac{1}{0} = \infty \\
 JS(\mathbb{P}||\mathbb{Q}) &= \frac{1}{2} KL\left(\mathbb{P}||\frac{\mathbb{P}+\mathbb{Q}}{2}\right) + \frac{1}{2} KL\left(\mathbb{Q}||\frac{\mathbb{Q}+\mathbb{P}}{2}\right) \\
 &= \frac{1}{2} \left( \sum_{x=0, Y \sim U[0,1]} 1 \times \ln \frac{1}{0.5} + \sum_{x=\theta, Y \sim U[0,1]} 1 \times \ln \frac{1}{0.5} \right) \\
 &= \ln 2 \\
 W(\mathbb{P}, \mathbb{Q}) &= |\theta|
 \end{aligned}$$

If  $\theta = 0$ , then  $\mathbb{P}$  and  $\mathbb{Q}$  are completely overlapped. This indicates that in most cases,  $KL$  divergence is impossible to calculate when two distributions are discrete.

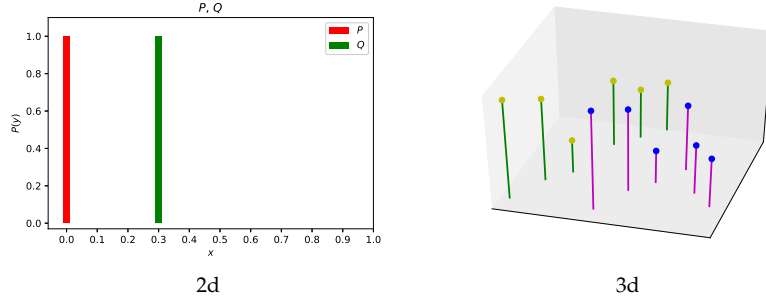


Figure 1: 2d and 3d plot of non-overlapped probability distributions  $\mathbb{P}$  and  $\mathbb{Q}$

#### Lebesgue's Dominated Convergence Theorem.

Let  $\{f_n\}$  be a sequence of real-valued measurable functions on a measure space  $(S, \Sigma, \mu)$ . Suppose that the sequence converges pointwise to a function  $f$

$$|f_n(x)| \leq f(x)$$

for all the index number  $n \in \mathbb{N}^+$  and  $x \in S$ , then  $f$  is integrable and

$$\lim_{n \rightarrow \infty} \int_S |f_n - f| d\mu = 0$$

#### Radamacher's theorem

If  $U$  is an open subset in  $\mathbb{R}^n$ ,  $f : U \mapsto \mathbb{R}^m$  is Lipschitz continuous, then  $f$  is differentiable almost everywhere.

#### Assumption 1

Let  $g : \mathcal{Z} \times \mathbb{R}^d \mapsto \mathcal{X}$  be locally Lipschitz between finite dimensional vector spaces,  $g_\theta(z)$  is denoted as evaluation on coordinates  $(z, \theta)$ . We call  $g$  satisfies assumption 1 for a certain probability distribution  $\mathbb{P}$  over  $\mathcal{Z}$  if there are local Lipschitz constants  $K(\theta, z)$  such that:

$$\mathbb{E}_{z \sim \mathbb{P}} [K(\theta, z)] < \infty$$

#### Theorem 1

Let  $\mathbb{P}_r$  be a fixed distribution over  $\mathcal{X}$ . Let  $Z$  be a random variable (e.g Gaussian) over another space  $\mathcal{Z}$ . Let  $g : \mathcal{Z} \times \mathbb{R}^d \mapsto \mathcal{X}$  be a function, that will be denoted  $g_\theta(z)$  with  $z$  the first coordinate and  $\theta$  the second. Let  $\mathbb{P}_\theta$  denote the distribution of  $g_\theta(Z)$ . Then:

1. If  $g$  is continuous in  $\theta$ , then  $W(\mathbb{P}_r, \mathbb{P}_\theta)$  is also continuous.
2. If  $g$  is locally Lipschitz and satisfies regularity assumption 1, then  $W(\mathbb{P}_r, \mathbb{P}_\theta)$  is continuous everywhere, and differentiable almost everywhere.

Proof.

1.

The goal is to show

$$\begin{aligned} & \lim_{\theta_1 \rightarrow \theta_2} \|g_{\theta_1}(z) - g_{\theta_2}(z)\| = 0 \\ \implies & \lim_{\theta_1 \rightarrow \theta_2} \|W(\mathbb{P}_r, \mathbb{P}_{\theta_1}) - W(\mathbb{P}_r, \mathbb{P}_{\theta_2})\| = 0 \end{aligned}$$

Assume  $\theta_1$  and  $\theta_2$  are two parameter vectors in  $\mathbb{R}^d$ . By the definition of the Wasserstein distance,

$$\begin{aligned} W(\mathbb{P}_{\theta_1}, \mathbb{P}_{\theta_2}) & \leq \mathbb{E}_{(a,b) \sim \gamma} [\|x - y\|] \\ & = \mathbb{E}_{(a,b) \sim \gamma} [\|g_{\theta_1}(z) - g_{\theta_2}(z)\|] \\ & = \int_{\mathcal{X} \times \mathcal{X}} [\|g_{\theta_1}(z) - g_{\theta_2}(z)\|] d\gamma \end{aligned} \tag{1}$$

If  $g$  is continuous, by the property of continuity,

$$\begin{aligned} & \lim_{\theta_1 \rightarrow \theta_2} g_{\theta_1}(z) = g_{\theta_2}(z) \\ \implies & \|g_{\theta_1}(z) - g_{\theta_2}(z)\| = 0. \end{aligned}$$

Since we assumed that  $\mathcal{X}$  is compact, which means closed and bounded, therefore exists a constant  $M$  such that  $\|g_{\theta_1}(z) - g_{\theta_2}(z)\| \leq M$  uniformly for all  $\theta \in \mathbb{R}^d$  and  $z$ . By the Lebesgue's dominated convergence theorem and (1),

$$\lim_{\theta_1 \rightarrow \theta_2} \int_{\mathcal{X} \times \mathcal{X}} [\|g_{\theta_1}(z) - g_{\theta_2}(z)\|] d\gamma = \lim_{\theta_1 \rightarrow \theta_2} \mathbb{E}_z [\|g_{\theta_1}(z) - g_{\theta_2}(z)\|] = 0$$

by the triangle inequality:

$$\|W(\mathbb{P}_r, \mathbb{P}_{\theta_1}) - W(\mathbb{P}_r, \mathbb{P}_{\theta_2})\| \leq W(\mathbb{P}_{\theta_1}, \mathbb{P}_{\theta_2})$$

Therefore if  $\theta_1 \rightarrow \theta_2$ ,

$$\lim_{\theta_1 \rightarrow \theta_2} W(\mathbb{P}_{\theta_1}, \mathbb{P}_{\theta_2}) = 0$$

The continuity of  $W(\mathbb{P}_r, \mathbb{P}_\theta)$  is proved.

For proving 2, assumed  $g$  is Lipschitz locally, that is, for a given point  $(z, \theta)$  in the domain of  $g$ , exists an open set  $(z, \theta) \in U$  and a constant  $K(z, \theta)$  such that  $\forall (z', \theta') \in U$  satisfied the property,

$$\|g_\theta(z) - g_{\theta'}(z')\| \leq K(z, \theta) \|(\theta, z) - (\theta', z')\| \leq K(z, \theta) (\|\theta - \theta'\| + \|z - z'\|).$$

After taking expectation for both sides of the above equation and set  $z = z'$ ,

$$\mathbb{E}_{z \sim \mathbb{P}_\theta} [\|g_\theta(z) - g_{\theta'}(z)\|] \leq \mathbb{E}_{z \sim \mathbb{P}_\theta} [K(z, \theta)] \|\theta - \theta'\| \quad \forall (z, \theta') \in U$$

define a new set

$$U_\theta := \{\theta' \mid (\theta', z) \in U\}.$$

Since  $U$  is open set, thus  $U_\theta$  as well. Based on the above proof, by assumption 1 and define  $K(\theta) = \mathbb{E}_{z \sim \mathbb{P}_\theta} [K(z, \theta)]$  we have:

$$\begin{aligned} \|W(\mathbb{P}_r, \mathbb{P}_{\theta_1}) - W(\mathbb{P}_r, \mathbb{P}_{\theta_2})\| &\leq W(\mathbb{P}_{\theta_1}, \mathbb{P}_{\theta_2}) \\ &\leq \mathbb{E}_{z \sim \mathbb{P}_\theta} [\|g_\theta(z) - g_{\theta'}(z')\|] \\ &\leq K(\theta) \|\theta - \theta'\| \quad \forall \theta' \in U_\theta \end{aligned}$$

this implies that  $W(\mathbb{P}_r, \mathbb{P}_{\theta_1})$  is locally Lipschitz, and also continuous everywhere. By Radamacher's theorem,  $W(\mathbb{P}_r, \mathbb{P}_{\theta_1})$  must be differentiable almost everywhere.  $\square$

There is a corollary guarantees that Wasserstain Distance can be used as loss function in neural networks.

### Corollary 1

If  $g_\theta$  is a feedforward neural network with  $\theta$  as its parameters, that is,  $g_\theta$  is a function composed by affine transformations and non-linear activation functions which are smooth Lipschitz continuous, and  $\mathbb{E}_{z \sim p(z)} [\|z\|] < \infty$ . Then the assumption 1 is satisfied and therefore  $W(\mathbb{P}_r, \mathbb{P}_\theta)$  is continuous everywhere and differentiable almost everywhere.

### 3.3 Wasserstein GAN

The solution of Wasserstein distance' infimum form in (1) is hard to find. However, solving (1) is equivalent to solve the Kantorovich - Rubinstein duality (Villani. 2008)

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r(\mathbf{x})} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_\theta(\mathbf{x})} [f(\mathbf{x})] \quad (5)$$

where all functions of  $f : \mathcal{X} \mapsto \mathbb{R}$  are 1-Lipschitz. If substitute  $\|f\|_L \leq 1$  to  $\|f\|_L \leq K$ , then the left side of the above equations would be  $K \cdot W(\mathbb{P}_r, \mathbb{P}_\theta)$ , and we are still in the same optimization problem. Thus, we could solving the Wasserstein distance by maximizing

$$\max_{w \in \mathcal{W}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r(\mathbf{x})} [f_w(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_\theta(\mathbf{x})} [f_w(g(\mathbf{z}))]$$

Where  $w$  are parameters for family of functions  $\{f_w\}_{w \in \mathcal{W}}$ , and  $w$  belong to a compact space  $\mathcal{W}$ , this implies that all  $f_w$  are  $K$ - Lipschitz functions depend only on  $\mathcal{W}$ . In order to implement that, Arjovsky et al. proposed to regulate all weights  $w$  of discriminator's neural network, e.g.,  $w \in [-0.01, 0.01]$  after update all weights. However, Arjovsky et al. also pointed out that clip weight is a bad way to implement Lipschitz constraint. Since if the clipping bound is too large, it would make lot of extra iterations for all weights  $w$  to converge their limit. On the contrary, if the clipping bound is too small, it would cause vanishing gradients. But still, since the Wasserstein distance is differentiable almost everywhere, ideally, the critic could be trained until optimality.

In addition, the backpropagation of  $G$  describe in the theorem below

**Theorem 3** [ Theorem 3 (Arjovsky et al. 2017) ]

Assume  $\mathbb{P}_r$  is a distribution, and  $\mathbb{P}_g$  is another distribution over  $g_\theta(Z)$  with  $Z$  as random variable,  $g_\theta(Z)$  satisfied assumption 1, then there exists a solution  $f : \mathcal{X} \mapsto \mathbb{R}$  for (5). Moreover,

$$\nabla_\theta W(\mathbb{P}_r, \mathbb{P}_\theta) = -\mathbb{E}_{\mathbf{z} \sim p_g(\mathbf{z})} [\nabla_\theta f(g_\theta(\mathbf{z}))].$$

With the above theoretical exploration, the final WGAN algorithm they present is like below

---

**Algorithm 2:** WGAN algorithm with default hyperparameters proposed by [Arjovsky et al.](#) Both  $D$  and  $G$  are use RMSprop as optimizer, with learn rate  $\lambda = 0.00005$ , weight clip  $c = 0.01$ ,  $n = 5$ , batch size  $k = 64$

---

```

1: while before  $\theta$  converged do
2:   for training discriminator  $n$  iterations do
3:     Sample  $k$  minibatch random noise from latent space distribution  $p_g(z)$ 
4:     Sample  $k$  minibatch input examples from data distribution  $p_r(x)$ 
5:     Backpropagating all weights in  $D$  by maximizing and update the loss:

```

$$g_{w_D} \leftarrow \nabla_{w_D} \frac{1}{k} \sum_{i=1}^k [f_w(x^{(i)}) - f_w(g(z^{(i)}))].$$

```

6:      $w_D \leftarrow w_D + \lambda \cdot \text{RMSProp}(w_D, g_{w_D})$ 
7:      $w_D \leftarrow \text{clip}(w_D, [-c, c]) \forall w_D \notin [-c, c]$ 
8:   end for
9:   Sample  $k$  minibatch random noise from latent space distribution  $p_g(z)$ 
10:  Backpropagating all weights in  $G$  by maximizing and update the loss:

```

$$g_{w_G} \leftarrow -\nabla_{w_G} \frac{1}{k} \sum_{i=1}^k [f_w(g(z^{(i)}))]$$

```

11:   $w_G \leftarrow w_G - \lambda \cdot \text{RMSProp}(w_G, g_{w_G})$ 
12: end while

```

---

The loss function is an approximation of Wasserstein distance



## Chapter 4

# Several improved GANs

### 4.1 Wasserstein GAN with Gradient Penalty (WGAN-GP)

In the previous chapter, we know that in Wasserstein GAN, [Arjovsky et al](#) introduced the new loss functions for  $D$  and  $G$  respectively

- $L(D) = -\mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_g(\mathbf{z})} [D(G(\mathbf{z}))]$
- $L(G) = \mathbb{E}_{\mathbf{z} \sim p_g(\mathbf{z})} [D(G(\mathbf{z}))]$

By the Kantorovich-Rubinstein duality, the output of discriminator much be enforced to a K-Lipschitz constraint. In other words, the L2 norm of  $D$ 's gradient must be bounded by a constant  $K$

$$\|\nabla_{\mathbf{x}} D(\mathbf{x})\| \leq K$$

In the previous chapter, [Arjovsky et al](#) introduced weight clipping

$$w \leftarrow \text{clip}(w, [-c, c])$$

to implement Lipschitz constrain, where  $c$  is the hyperparameter. However, as [Gulrajani et al](#) point out in their paper, this method have two problems:

1. gradient can easily vanish or explode
2. most of weights are distributed on the clipping boundary, i.e., either equals to  $-c$  or  $c$ .

In the first problem, since discriminator usually are deep neural networks, therefore, by the chain rule, if clipping constant  $c$  is being set relatively small, the backpropagation can cause gradient vanish; if  $c$  being set relatively big, gradient can easily explode. Instead of applying weight clipping, [Gulrajani et al](#) showed on Swiss Roll dataset (Figure 10) that the gradient penalty term can make gradient stably through each layer.

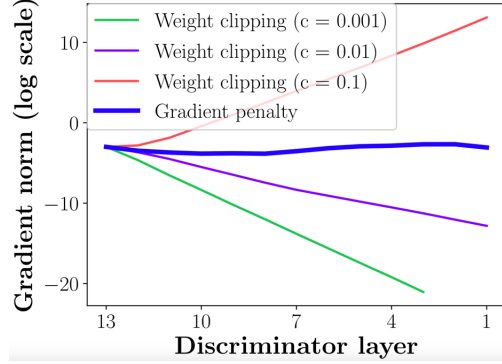


Figure 10:  $x$  axis indicate layer's index . Figure 1(b) in [Gulrajani et al \[18\]](#)

In the second problem, weight are bounded on a interval  $[-c, c]$ , i.e., for those weights greater than  $c$  will be set to  $c$ , for less than  $-c$  will set back to  $-c$ . This would cause discriminator to learn a simple function. [Gulrajani et al](#) verified the problem on Swiss Roll dataset and found that most of weights are distributed either  $c$  or  $-c$ . Instead, the gradient penalty term dose not suffer from this issue, Figure 11 demonstrate this.

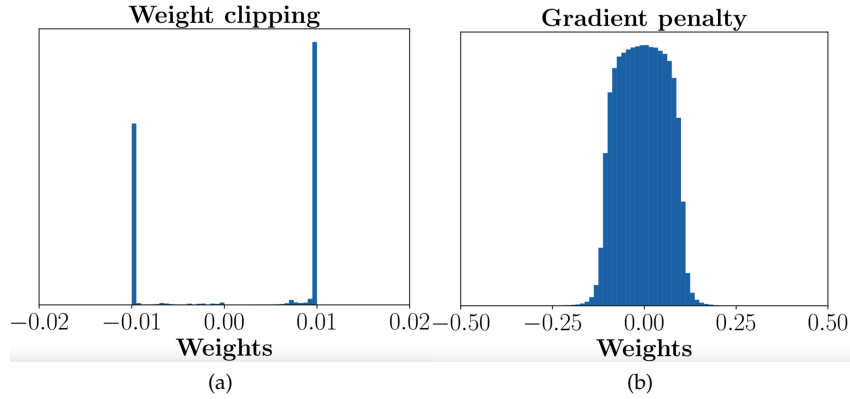


Figure 11: Figure 1(b) in [Gulrajani et al \[18\]](#)

The above two problem indicate that for the weight clipping method, model performance is too sensitive with respect to the corresponding hypermarameter  $c$ .

A differentiable function is  $K$ -Lipschitz if and only if the Euclidean norm of it's gradient is almost  $K$  everywhere. Instead of using weight clip, [Gulrajani et al](#) proposed to add a gradient penalty term which mentioned above to enforce 1-Lipschitz constraint.

$$(\|\nabla_x D(x)\| - 1)^2$$

## 4.2. Wasserstein GAN with A Consistency Term (CT-GAN)

when discriminator sufficiently trained, the norm of it's gradient will close to 1. Then the new loss function to be minimize becomes

$$L = \underbrace{\mathbb{E}_{z \sim \mathbb{P}_g} [D(g(z))] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)]}_{\text{original WGAN's loss}} + \underbrace{\lambda \cdot \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} \left[ \left( \|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1 \right)^2 \right]}_{\text{new gradient penalty term}}.$$

Where  $\hat{x} = \epsilon x + (1 - \epsilon)g(z)$  and  $\epsilon$  is a random real number from  $[0,1]$ . Thus model can be penalized by the gradient penalty term penalizes if the norm of gradient away from 1.

For experiment setting, [Gulrajani et al](#) proposed to not use batch normalization, but layer normalization is recommend. Below is the final algorithm:

---

**Algorithm 4:** [Algorithm 1 in [Gulrajani et al](#) [18] ]. The default hyperparameters are  $\lambda = 10, n = 5, \beta_1 = 0, \beta_2 = 0.9, \alpha = 0.0001$ , where  $\beta_1, \beta_2$  and  $\alpha$  are Adam optimizer hyperparameters

---

```

1: while before  $\theta$  converged do
2:   for training discriminator  $n$  iterations do
3:     for  $i = 1, \dots$ , batch size  $m$  do
4:       Sample  $x$  from data distribution  $p_r(x)$ 
5:       Sample  $z$  from latent space  $p_g(z)$ 
6:       Sample random number  $\epsilon$  from  $U[0, 1]$ 
7:        $\hat{x} \leftarrow \epsilon x + (1 - \epsilon)G_\theta(x)$ 
8:        $L^{(i)} \leftarrow D(G(z)) - D(z) + \lambda(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2$ 
9:     end for
10:     $w_D \leftarrow w_D - \eta \cdot \text{Adam}\left(\frac{1}{m} \sum_{i=1}^m \nabla_{w_D} L^{(i)}, \beta_1, \beta_2, \alpha\right)$ 
11:   end for
12:   Sample  $k$  minibatch random noise from latent space distribution  $p_g(z)$ 
11:    $w_G \leftarrow w_G - \eta \cdot \text{Adam}\left(-\frac{1}{m} \sum_{i=1}^m \nabla_{w_G} D(G(z^{(i)})), \beta_1, \beta_2, \alpha\right)$ 
12: end while

```

---

## 4.2 Wasserstein GAN with A Consistency Term (CT-GAN)

Despite the WGAN with gradient penalty term overcome the two problems which exist in the original WGAN, [Wei et al. \(2017\)](#) indicate that the gradient penalty term only works on those interpolated samples  $\hat{x}$ . At the early training stage, the distribution  $\mathbb{P}_g$  is far away from  $\mathbb{P}_r$ . Moreover, the gradient penalty term usually fail to exam the continuity near real sample  $x$ . Therefore discriminator can easily out of 1-Lipschitz continuity.

Based on these issues, they introduced to improve the WGAN with gradient penalty by enforce the Lipschitz continuity over the real data  $x \sim \mathbb{P}_r$ . Specifically, they perturb each real sample  $x$  as  $x'$  and  $x''$ , then apply Lipschitz constraint to bound the  $D(x')$  and  $D(x'')$ .

Recall from the WGAN with gradient penalty, a discriminator  $D : \mathcal{X} \mapsto \mathcal{Y}$

---

## 4.2. Wasserstein GAN with A Consistency Term (CT-GAN)

---

is Lipschitz continuous if and only if there exists a real number  $K$  such that  $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ ,

$$\|D(\mathbf{x}_1) - D(\mathbf{x}_2)\|_2 \leq K \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|_2$$

therefore, the following soft consistency term (CT) can be added in order to penalize those violations from the above inequality.

$$CT|_{\mathbf{x}_1, \mathbf{x}_2} = \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} \left[ \max \left( \frac{\|D(\mathbf{x}_1) - D(\mathbf{x}_2)\|_2}{\|\mathbf{x}_1 - \mathbf{x}_2\|_2} - M', 0 \right) \right]$$

However, [Wei et al](#) point out that  $\|\mathbf{x}_1 - \mathbf{x}_2\|$  is impractical, as it is hard to compute all the combinations of  $(\mathbf{x}_1, \mathbf{x}_2)$ . Therefore, to make the inequality can be implement, all  $(\mathbf{x}_1, \mathbf{x}_2)$  assumed be bounded by a constant  $M'$ , and  $M'$  is a hyper-parameter. In addition, they proposed to apply dropout before the output of discriminator, denoted as  $D(\mathbf{x}')$ , where  $\mathbf{x}'$  is sampled from input data, and  $D(\mathbf{x}'')$  is the second data point which applied in the same corresponding way.  $D_{-}(\cdot)$  denote the output of the second last layer of the discriminator. Thus the final consistency term and loss function is

$$CT|_{\mathbf{x}', \mathbf{x}''} = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} \left[ \max(0, \|D(\mathbf{x}_1) - D(\mathbf{x}_2)\|_2) + 0.1 \cdot \|(D_{-}(\mathbf{x}'), D_{-}(\mathbf{x}''))\|_2 - M' \right],$$

$$L = \mathbb{E}_{\mathbf{z} \sim \mathbb{P}_g} [D(G(\mathbf{z}))] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] + \lambda_1 GP|_{\hat{\mathbf{x}}} + \lambda_2 CT|_{\mathbf{x}', \mathbf{x}''}.$$

Final algorithm for the CT-GAN

---

**Algorithm 5:** [Algorithm 1 in [Wei et al](#) [15] ]. The default hyperparameters are  $\lambda_1 = 10$ ,  $\lambda_2 = 2$ ,  $n = 5$ ,  $\beta_1 = 0.5$ ,  $\beta_2 = 0.9$ ,  $\alpha = 0.0002$ ,  $m = 64$  where  $\beta_1, \beta_2$  and  $\alpha$  are Adam optimizer hyperparameters.

---

```

1: while before  $\theta$  converged do
2:   for training discriminator  $n$  iterations do
3:     for  $i = 1, \dots$ , batch size  $m$  do
4:       Sample  $\mathbf{x}$  from data distribution  $p_r(\mathbf{x})$ 
5:       Sample  $\mathbf{z}$  from latent space  $p_g(\mathbf{z})$ 
6:       Sample random number  $\epsilon$  from  $U[0, 1]$ 
7:        $\hat{\mathbf{x}} \leftarrow \epsilon \mathbf{x} + (1 - \epsilon) G_\theta(\mathbf{x})$ 
8:        $L^{(i)} \leftarrow D(G(\mathbf{z})) - D(\mathbf{z}) + \lambda_1 GP|_{\hat{\mathbf{x}}} + \lambda_2 CT|_{\mathbf{x}', \mathbf{x}''}$ 
9:     end for
10:     $w_D \leftarrow w_D - \eta \cdot \text{Adam} \left( \frac{1}{m} \sum_{i=1}^m \nabla_{w_D} L^{(i)}, \beta_1, \beta_2, \alpha \right)$ 
11:   end for
12:   Sample  $k$  minibatch random noise from latent space distribution  $p_g(\mathbf{z})$ 
11:    $w_G \leftarrow w_G - \eta \cdot \text{Adam} \left( -\frac{1}{m} \sum_{i=1}^m \nabla_{w_G} D(G(\mathbf{z}^{(i)})), \beta_1, \beta_2, \alpha \right)$ 
12: end while

```

---

## Chapter 5

# Experiments

[Conneau et al](#) showed results trained by Vanilla GAN for some language pairs such as English(en) translate to Spanish(es), English to French(fr) and English to Esperanto(eo) etc. Their code, dictionaries and word embeddings are open source on GitHub ([MUSE](#))<sup>1</sup>.

Since the GPU computational resource is limited, and for the convenient of comparing the result and benchmark from the MUSE, we only select several relatively representative language pairs which [Conneau et al. \(2017\)](#) implemented, and plus few other language pairs.

We will describe the data we used and review the experimental procedure in details from [Conneau et al](#) first, then show our result.

### 5.1 data description

Facebook’s AI Research department published pre-trained word vectors for 294 languages in GitHub. They trained them on Wikipedia as source text using the library called fastText. More specifically, the word vectors that were trained using the skip-gram models ([Bojanowski et al. 2016](#)) with 300 as embedding dimension. All dataset can be download from their GitHub repository<sup>2</sup>.

We also tried to fit English embeddings which are trained from different corpus, this target embeddings are trained from English Common Crawl Corpus (ENC3)<sup>3</sup>, which contains two million word embeddings.

### 5.2 Multilingual Unsupervised and Supervised Embeddings

In the unsupervised method of MUSE, the linear transformation  $W$  across source language embeddings and target embeddings can be learned by GANs without any given dictionary. The experimental setup will be review below.

---

<sup>1</sup>[MUSE: Multilingual Unsupervised and Supervised Embeddings](#)

<sup>2</sup>[FastText: Pre-trained word vectors](#)

<sup>3</sup>[NLPL word embeddings repository](#)

### 5.2.1 Network architecture

The vanilla GAN's neural architecture [Conneau et al. \(2017\)](#) proposed is to define the generator as a linear transformation  $W$  which maps source embedding matrix  $X$  to target embedding matrix  $Y$ , therefore the generator is a single-layer neural network without activation function, the number of nodes in that layer depends on the embedding dimension of language pairs. Since the dimensionality of all monolingual embeddings is 300, hence the generator is a matrix with size  $300 \times 300$ , it is a square matrix. The discriminator is a multilayer perceptron with two hidden layers, both have 2048 nodes with Leaky-ReLu as activation function. The output of discriminator is single node squeezed to probability by sigmoid.

### 5.2.2 Objective functions

Assume  $\mathcal{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$  and  $\mathcal{Y} = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m)}\}$  are source and target language embeddings respectively, the label for source embedding is 1, 0 for target. Since the goal of the discriminator is to identify whether the input coming from source or target. Therefore the discriminator's loss function is

$$\mathcal{L}_D = -\frac{1}{n} \sum_{i=1}^n \log P_{\theta_D}(\text{source} = 1 | W\mathbf{x}^{(i)}) - \frac{1}{m} \sum_{i=1}^m \log P_{\theta_D}(\text{target} = 0 | \mathbf{y}^{(i)})$$

where  $\theta_D$  denote the parameters of discriminator's network,  $P_{\theta_D}(\text{source} = 1 | W\mathbf{x})$  means probability of input point coming from source embeddings, whereas  $P_{\theta_D}(\text{target} = 0 | \mathbf{y})$  means target embeddings, i.e., real dataset.

The task of generator is to deceive the discriminator as much as possible, thus, it's loss function is

$$\mathcal{L}_W = -\frac{1}{n} \sum_{i=1}^n \log P_{\theta_D}(\text{source} = 0 | W\mathbf{x}^{(i)}) - \frac{1}{m} \sum_{i=1}^m \log P_{\theta_D}(\text{target} = 1 | \mathbf{y}^{(i)})$$

### 5.2.3 Optimizer

Since they use minibatch for each training iteration and the batch size is 32, stochastic gradient descent is applied for minimizing  $\mathcal{L}_D$  and  $\mathcal{L}_W$ , the default learning rate is 0.1 without momentum.

### 5.2.4 Orthogonality

Since the generator is a  $300 \times 300$  matrix, [Conneau et al. \(2017\)](#) proposed to applied an orthogonal constraint to generator after each minibatch iteration. It keep embeddings' quality after source embedding be translated. Moreover, orthogonal matrix have a nice property, assume  $\mathcal{O}$  is an orthogonal matrix, i.e.  $\mathcal{O}^T \mathcal{O} = \mathcal{O} \mathcal{O}^T = I$ , it preserve  $\ell_2$  norm,

$$\|\mathcal{O}\mathbf{x}\|_2^2 = \langle \mathcal{O}\mathbf{x}, \mathcal{O}\mathbf{x} \rangle = \mathcal{O}^2 \langle \mathbf{x}, \mathbf{x} \rangle = \|\mathbf{x}\|_2^2$$

This implies  $\|\mathcal{O}x\|_2 = \|x\|_2$ , in addition this means  $\mathcal{O}$  is an isometry in  $\ell_2$  space because it preserve distance.

The algorithm they used to update transformation  $W$  is

$$W \leftarrow (1 + \beta)W - \beta(WW^T)W$$

where  $\beta$  is a hyperparameter, as they suggest,  $\beta = 0.01$  performs well.

### 5.2.5 Cross-domain similarity local scaling (CSLS)

A reliable metric is crucial for the criterion of model selection, as we need to compare the similarity between source embeddings and target embeddings.

Nearest neighbors have the asymmetric property, which means that if  $x$  is the  $k$ -nn of a point  $y$ , does not imply that  $y$  is a  $k$ -nn of  $x$ . By [Radovanović et al. \(2010\)](#), the asymmetric property would cause harmful result in high dimensional spaces: For some points, which are called hubs, are nearest neighbors of numerous of other points, however, for some other points, which are called anti-hubs, are even no any points are nearest neighbors of them.

Therefore, [Conneau et al. \(2017\)](#) proposed a bipartite neighborhood graph called Cross-domain similarity local scaling(CSLS) as the new metric. By their definition, the mean cosine similarity from a translated source embedding  $Wx_s$  to its  $k$ -target neighborhood in target space is defined as

$$r_T(Wx_s) = \frac{1}{k} \sum_{y_t \in \mathcal{N}_T(Wx_s)} \cos(Wx_s, y_t)$$

where  $\mathcal{N}_T(Wx_s)$  denote the neighborhood associated with a translated source word  $Wx_s$ , and all  $k$  words of  $\mathcal{N}_T(Wx_s)$  are from target embeddings. Similarly,  $\mathcal{N}_S(y_t)$  denote the neighborhood in target embedding corresponded to a target word  $y_t$ . The mean cosine similarity of a target word to its  $k$ -target neighborhood in translated source embeddings space is defined as

$$r_S(y_t) = \frac{1}{k} \sum_{Wx_s \in \mathcal{N}_S(y_t)} \cos(y_t, Wx_s)$$

these two mean cosine scores for translated source embeddings and target embeddings are implemented by the library called Faiss<sup>4</sup>, which greatly speed up the computation of nearest neighbors. The Faiss is developed by [Johnson et al. \(2017\)](#) who are members of Facebook AI research group. The CSLS algorithm define as

$$\text{CSLS}(Wx_s, y_t) = 2\cos(Wx_s, y_t) - r_T(Wx_s) - r_S(y_t).$$

Obviously, there is no any hyperparameter tuning requirement for it. By [Conneau et al. \(2017\)](#) interpretation, the update of CSLS algorithm increase the similarity for those anti-hub word embeddings. In contrast, it also reduce those embeddings in dense hub regions. In addition, the CSLS greatly improves the translation accuracy in experiments.

<sup>4</sup>Faiss: a package for efficient similarity computation.

### 5.2.6 Translation process

The MUSE will translate source embeddings to the target embedding space, then for each translated source embedding, select the k-nearest neighbors from target embeddings, by CSLS algorithm.

### 5.2.7 Evaluation procedure

Conneau et al provides some bilingual evaluation dictionaries which are translated by an internal translation tool, each dictionary translate 1500 unique source words, most of them are commonly used, also polysemy for some vocabularies, i.e., the English word "discover" translate to the Spanish words as "descubrir", "descubra", "descubre" and "descubren". After each epoch, the MUSE will calculate the CSLS accuracy on bilingual evaluation dictionary, and report the precision at 1, 5 and 10, i.e., the accuracy of nearest neighbors.

MUSE will also compute the similarity for the top 10,000 frequent words. After finished the calculation of CSLS precision, the length of both source and target embeddings will be normalized to 1, then a subset of source word index which correspond to the temporary dictionary will be selected (source and target word embeddings would be indexed as temporary dictionary before adversarial training), and transformed by learned generator, to match the target word by k-nearest neighbor and CSLS algorithm. At the end of the epoch, the best model will be save based on the mean cosine csls for the maximum of 10000 top frequent embeddings.

### 5.2.8 Procruster refinement

Use those most frequent words as anchor points, then apply Procruster algorithm, i.e., compute the  $SVD(YX^T)$ , where  $X$  and  $Y$  is source and target embeddings respectively, then use  $U$  and  $V$  to solve

$$W^* = \underset{W}{\operatorname{argmin}} \|WX - Y\|_F = UV^T,$$

## 5.3 Experimental set up

In our experiment, we focused on three improved GANs, i.e., Wasserstein GAN, WGAN with gradient penalty and CT-GAN, to compare the result from vanilla GAN. Our code is based on MUSE, in which use Pytorch as deep learning library instead of TensorFlow. We later found that for WGAN, add layer normalization for the output of fully connected layer can not only increase the accuracy, but also stabilize the training process. Therefore, we would discuss with layer normalization and without layer normalization for these three improved GANs later.

### 5.3.1 Target language selection

We fix English as the source language, and select all the languages from the



table 2 of Søgaard et al (2018) as target languages, which are Spanish(ES) Estonian(ET), Greek(EL), Finnish(FI), Hungary(HU)and Polish(PL). Plus Chinese(ZH), which belongs to a different language family, and English which the embeddings are trained from different Corpus.

### 5.3.2 Hyper-parameter search

Since the hyper-parameter search cost lot of GPU resource, therefore, we decide to only to tune those main hyper-parameters which affect most for our model, for those secondary hyper-parameters such as  $\beta_1, \beta_2$  and  $\alpha$  in Adam optimizer, we applied the default from PyTorch. We first use grid search to find relatively good hyper-parameters for English to Spanish, Greek and Estonian for all three different Algorithms, then extends to other languages pairs. For Wasserstein GAN, we found  $c = 0.25$ , and RMSprop with learning rate 0.0005 works well.

For WGAN with gradient penalty, we use  $\lambda = 0.5$ , and Adam with learning rate 0.0005.

For CT-GAN, we use  $\lambda_{GP} = 0.5$ ,  $\lambda_{CT} = 2$  and  $m = 0.2$ .

### 5.3.3 Define the metric

We use two types of metrics which mentioned at the previous section to select models and hyper-parameters.

- The CSLS nearest neighbor precision at 1, 5 and 10, which evaluate the performance on the evaluation dictionary.
- The mean cosine CSLS, which measure the similarity of the top frequent 10,000 words.

In table 5.1, the source language is English, we can see that the mean cosine CSLS inaccurately reflect translation performance, e.g. for English to Chinese(zh) and Estonian(et), the precision 10 is very low, but has high overall mean CSLS, this indicate that the translated embeddings have very close neighbors, but almost all of them are incorrect translations. Therefore, we use CSLS nearest neighbor (P@1,5 and 10) to evaluate translation performance.

Vanilla	en	es	et	el	fi	hu	tr	pl	zh
p@1	0.0	39.8	3.2	13.9	10.3	19.1	16.4	20.3	0.0
p@5	0.0	67.8	9.8	28.4	26.3	38.6	31.8	46.5	0.0
p@10	0.0	75.1	14.2	36.4	35.8	48.1	39.9	57.4	0.0
mean csls	47.9	67.1	49.6	56	56.1	56.8	55.5	58.6	79.6

Table 5.1: p@5, p@10 and mean cosine CSLS for differen target languages

### 5.3.4 Network architecture

We keep to used the same network architecture for all three improved GANs, i.e., two hidden layers with 2048 dimensionality, but, different loss functions from Vanilla GAN, and without sigmoid to squeeze discriminator's output to probability. Plus added layer normalization to each layers.

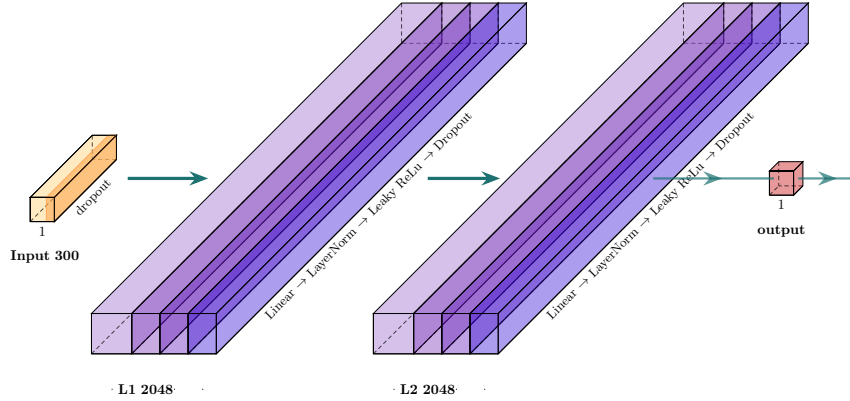


Figure 5.1: Architecture of discriminator

### 5.3.5 Stability test

Since the performance depends on how the weights are initialized, e.g., an experiment can received a good result, but in the same setting with different random weight initialization, the generator may failed. Therefore, for some algorithms, we fixed hyper-parameters and varying ten different random seeds to test how stable these GANs are. In addition, we define a failed running as if the corresponding P@10 less than 0.1.

### 5.3.6 Benchmark

We tried to reproduced the result in table 1 of [Conneau et al](#) , but we could not get the same result. This may due to the best random seed they did not specified, hence we implement the same setting from their paper, and use the new reproduced results as our benchmark.

## 5.4 Results

Since we later found that added layer normalization in each of discriminator's layer can moderately increase the performance, therefore we first report the results of three improved GANs with layer normalization, and later discuss each improved GAN individually. We also found that instead of exponentially shrinking the learning rate for each epochs, fix the learning rate until finished the training process performs better.

WGAN-GP	en	es	et	el	fi	hu	tr	pl	zh
p@1	83.7	36.3	11.2	13.2	14.5	23	16.2	19.6	5.7
p@5	89.3	64.3	24.2	29.5	34.3	44	31.1	43.8	12.6
p@10	90.8	71.3	30.9	36.6	42.3	52.5	38.5	53.8	16

Table 5.2: p@1 and p@5, p@10 of WGAN with GP, random seed is -1

WGAN	en	es	et	el	fi	hu	tr	pl	zh
p@1	0.0	40.4	1	7.7	4.1	16.5	0.0	0.0	0.0
p@5	0.1	68.6	3.5	19.2	12.4	34.8	0.0	0.0	0.1
p@10	0.1	75.6	6.2	26.5	17	43.5	0.1	0.1	0.1

Table 5.3: p@1 and p@5, p@10 WGAN, random seed is -1

CT-GAN	en	es	et	el	fi	hu	tr	pl	zh
p@1	74.0	27.2	6.0	8.4	11.1	16.5	10.3	13.8	0.2
p@5	80.7	50.8	14.9	20.9	25.7	34	24.2	34.6	1.1
p@10	82.4	58.5	19.6	28.0	32.9	41.9	30.5	43.8	1.8

Table 5.4: p@1 and p@5, p@10 CT-GAN, random seed is -1

We can see that WGAN with gradient penalty(table 5.2) works best and quite stable, and the result of English translate to itself has very good performance, and even to Chinese, WGAN-GP still learned part of its distribution instead of completely failed. In the results of vanilla GAN (table 5.1), the training from en to en and zh are completely failed, even for en to et which is not a far away language, it still barely learned small distribution from it.

we can see that WGAN-GP is not easily suffered from different dataset, since for en to en and zh, vanilla GAN completely failed

the original WGAN (table 5.3) works not as good as we expect.

## 5.4.1 Vanilla GAN

## 5.4.2 stability test

Vanilla	en	es	et	el	fi	hu	tr	pl	zh
Failed	10	1	9	0	0	1	1	0	10
max p@1	0.1	39.6	5.7	13.7	13.5	25.6	17.2	24.9	0.0
min p@1	0.0	37.9	0.0	12.0	9.1	21.5	10.9	17.3	0.0
avg p@1	0.0	39.2	0.6	12.2	11.1	22.7	15.0	19.5	0.0
std p@1	0.0	0.7	1.7	1.0	1.6	1.2	2.1	2.6	0.0
max p@5	0.1	67.7	15.6	30.0	31.7	45.2	34.0	50.7	0.1
min p@5	0.1	63.8	0	25.1	23.4	39.6	24.5	40.1	0.0
avg p@5	0.1	66.8	1.6	27.4	27.5	42.4	30.1	43.4	0.0
std p@5	0.0	1.2	4.7	1.8	2.5	1.6	3.1	3.4	0.0
max p@10	0.4	74.9	22.4	37.7	65.7	54.5	0.103	61.0	0.1
min p@10	0.1	71.4	0	31.8	31.3	47.5	30.9	50.3	0.0
avg p@10	0.1	74.1	2.3	34.8	35.5	51.3	37.5	53.8	0.1
std p@10	0.1	1.1	6.7	2.0	2.9	1.8	3.7	3.5	0.0

Table 4: result of ten runs Vanilla GAN with differen random seed

WGAN-GP	en	es	et	el	fi	hu	tr	pl	zh
Failed	1	0	0	0	0	1	0	0	0
max p@1	86.9	35.5	11.0	14.2	16.0	23.7	15.4	18.6	6.9
min p@1	0.0	34.6	10.3	6.7	13.4	0.0	11.9	16.2	3.4
avg p@1	76.3	35.0	9.9	8.6	14.5	19.7	14.7	17.8	5.3
std p@1	25.6	0.4	1.1	2.0	1.0	6.7	1.4	0.7	1.2
max p@5	91.8	62.2	24.0	31.0	36.0	44.9	32.6	43.8	14.2
min p@5	0.1	59.7	21.0	17.9	31.5	0	27.4	40.9	9.1
avg p@5	80.9	61.3	22.5	20.7	33.6	38.1	30.6	42.6	11.9
std p@5	27.0	0.8	1.2	3.6	1.4	12.8	1.7	0.9	1.9
max p@10	93.6	69.7	31.7	38.3	44.5	53.5	40.9	54.5	18.6
min p@10	0.4	67.5	27.0	23.3	39.5	0.1	35.6	50.9	12.0
avg p@10	82.1	68.7	29.1	26.9	42.1	45.8	38.3	52.7	15.3
std p@10	27.3	0.6	1.3	4.0	1.5	15.3	1.5	1.2	2.1

Table 4: result of ten runs Vanilla GAN with differen random seed

CT-GAN	en	es	et	el	fi	hu	tr	pl	zh
Failed	7	0	0	0	1	0	0	lack	10
max p@1	71.9	32.1	5.8	10.5	10.9	15.7	10.8	4.5	0.1
min p@1	0.0	28.8	4.2	8.3	0.0	13.4	8.9	0.0	0.0
avg p@1	21.9	30.3	5.2	9.2	9.1	15.3	9.4	0.5	0.0
std p@1	32.8	1.2	0.7	1.1	3.1	0.9	0.9	1.3	0.0
max p@5	80.7	54.8	15.7	23.9	26.9	34.8	23.8	12.3	0.4
min p@5	0.0	51.3	12.5	19.2	0.1	30.4	20.2	0.0	0.0
avg p@5	24.7	53.3	13.9	21.8	22.2	33.2	22.0	1.3	0.2
std p@5	36.2	1.3	1.1	1.8	7.5	1.4	1.2	3.7	0.1
max p@10	83.2	63.0	22.0	30.4	35.6	44.6	30.5	18.0	0.7
min p@10	0.0	58.7	17.4	24.4	0.1	39.1	25.7	0.0	0.1
avg p@10	25.6	61.0	19.6	27.8	28.9	41.4	28.4	1.9	0.3
std p@10	36.8	1.3	1.4	2.0	9.7	1.6	1.5	1.5	0.1

Table 4: result of ten runs WGAN original for different random seed

### 5.4.3 WGAN-Original

WGAN	en	es	et	el	fi	hu	tr	pl	zh
p@1	0.0	20.4	0.1	0.0	0.0	0.0	0.0	10.5	0.0
p@5	0.1	38.6	0.1	0.0	0.0	0.0	0.0	26.9	0.0
p@10	0.3	44.7	0.2	0.1	0.1	0.0	0.1	35.1	0.0

Table 4: p@1 and p@5, p@10 from different target languages

WGAN 10 ori	en	es	et	el	fi	hu	tr	pl	zh
Failed	10	3	10	10	10	2	8	9	10
max p@1	0.3	19.9	0.0	0.0	0.0	8.5	3.1	4.5	0.0
min p@1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
avg p@1	0.0	14.0	0.0	0.0	0.0	6.6	0.8	0.5	0.0
std p@1	0.1	9.2	0.0	0.0	0.0	3.3	1.2	1.3	0.0
max p@5	0.9	39.2	0.0	0.0	0.0	22.0	10.1	12.3	0.0
min p@5	0.1	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0
avg p@5	0.2	26.7	0.0	0.0	0.0	16.5	2.6	1.3	0.0
std p@5	0.2	17.4	0.0	0.0	0.0	8.3	4.0	3.7	0.0
max p@10	1.2	47.2	0.1	0.0	0.0	29.1	14.2	18.0	0.0
min p@10	0.1	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0
avg p@10	0.4	31.7	0.0	0.0	0.0	22.1	3.8	1.9	0.0
std p@10	0.3	20.5	0.0	0.0	0.0	11.1	5.8	5.4	0.0

Table 4: result of ten runs WGAN original for different random seed

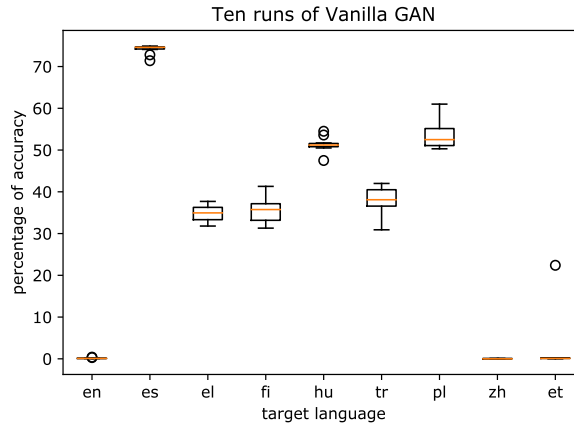


Figure 5.2: ten runs of vanilla GAN

**5.4.4 without shrinking learning rate****5.4.5 with layer normalization****5.4.6 Wasserstein GAN****5.4.7 result of other people did****5.4.8 comparison tables****5.4.9 Overview visualization**

T-Distributed Stochastic Neighbor Embedding (T-SNE) is a dimensionality reduction technique, which is a non-linear manifold learning algorithm. It developed by (Maaten et al. 2008). Compare with PCA, T-SNE more accurately captures main information from high dimensional space than PCA does. However, T-SNE costs much more computational time than PCA, hence, for visualize the performance of GANs worked on some languages pairs, we first extract the most 150 relevant components from 300 by PCA, then transform those best 150 components to 2d space by T-SNE.

We randomly sample 2000 points out of the most frequent 10,000 word embeddings from

- Original source embeddings before transformed (green dots).
- Source embeddings transformed by generator (blue dots).
- Target embeddings (red dots).

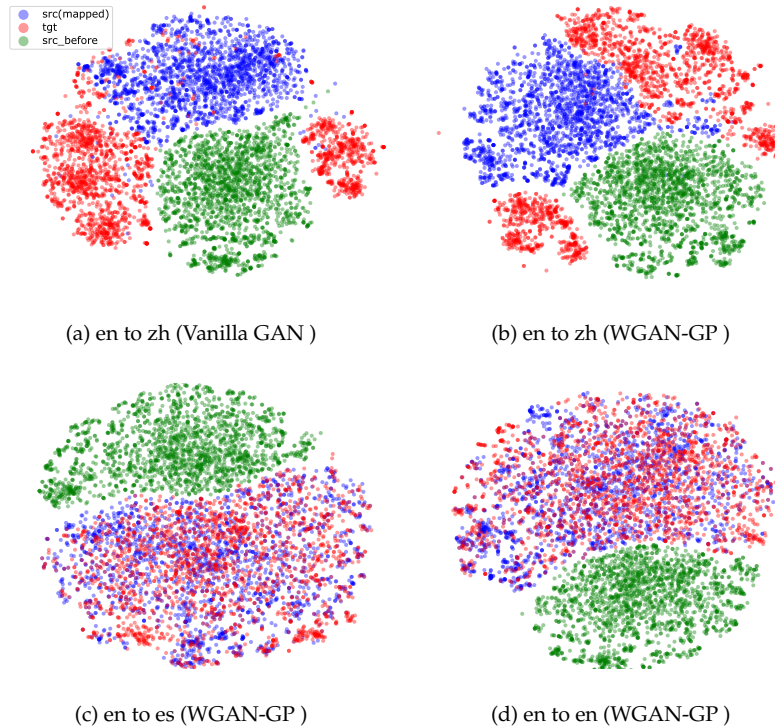


Figure 5.5: 2d T-SNE after extract 150 components by PCA

For English translate to relatively closed embeddings such as Spanish and English which is trained from different Corpus, we can see that from figures 5.5(a) and (b), the translated English embeddings are mixed very well with target embeddings, but for English translate to largely different language such as Chinese as showed in the figures 5.5(a) and (b), the discriminator successfully separate almost all samples, at least in the circumstance of linear generator.

#### 5.4.10 drawback of the three new WGANs

maybe there is no linear transformation between those languages, because EN to EN ,future work maybe remove the assumption, change the generator to non-linearity

#### 5.4.11 Final experiment

- translate some most frequent words,
- translate some rare words.
- compare how model performs on frequent words, and rare words.

## **Chapter 6**

# **Conclusions**

### **6.1 Future work**

investigate the non-linear transformation,



## Appendix A

# Some deep learning techniques

### A.1 Layer normalization

Training a deep neural network can cause internal covariate shift, i.e., the change of the input distribution of all layers. Especially for a deep neural network, even a small change can be expanded when the network goes deeper, therefore vanishing gradients can easily happen for saturating non-linearities. Batch normalization can help to reduce internal covariate shift by normalizing mini-batch input across their corresponding size. In particular,

$$\begin{aligned}\mu_j &= \frac{1}{m} \sum_{j=1}^n \mathbf{x}_{ij} \\ \sigma_j^2 &= \frac{1}{m} \sum_{j=1}^m (\mathbf{x}_{ij} - \mu_i)^2 \\ \hat{\mathbf{x}} &= \frac{\mathbf{x}_{ij} - \mu_i}{\sqrt{\sigma^2 + \epsilon}},\end{aligned}$$

where  $i$  indicates the index of example, and  $j$  is the index of feature. Thus the network converges relatively faster. The batch normalization method is proposed by [Lofte et al. \(2015\)](#).

Despite it achieving great success, batch normalization still has two main drawbacks:

- Model performance depends on the mini-batch size, because the mean and the variance, which are calculated by mini-batch, cause errors and these errors will vary between different mini-batches.
- It's very difficult to apply to RNN.

Layer normalization is proposed by [Ba et al. \(2016\)](#). Unlike batch normalization, layer normalization does not depend on the size of the mini-batch, and also works well on RNNs. For each input example, it normalizes across

features in a layer, more specifically,

$$\begin{aligned}\mu_i &= \frac{1}{m} \sum_{j=1}^m \mathbf{x}_{ij} \\ \sigma_i^2 &= \frac{1}{m} \sum_{j=1}^m (\mathbf{x}_{ij} - \mu_i)^2 \\ \hat{\mathbf{x}} &= \frac{\mathbf{x}_{ij} - \mu_i}{\sqrt{\sigma^2 + \epsilon}}\end{aligned}$$

## A.2 RMS prop

RMSprop is a optimization algorithm of gradient descent. It is proposed in the lecture 6 of [Hinton et al](#)<sup>5</sup> online course, and it is unpublished.

The motivation for developing RMSprop is to solve Adagrad's vanishing learning rate problem. Unlike other gradient descent algorithm such as sgd or mini-batch gradient descent, by [Ruder](#)'s explanation, Adagrad can adapt learning rate during training process, but the denominator term tends to zero after some training iterations. In particular, the RMSprop is

$$\begin{aligned}\mathbb{E}[\mathbf{g}^2]_t &= 0.9\mathbb{E}[\mathbf{g}^2]_{t-1} + 0.1\mathbf{g}_t^2 \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \frac{\eta}{\sqrt{\mathbb{E}[\mathbf{g}^2]_t + \epsilon}} \mathbf{g}_t\end{aligned}$$

where the  $\epsilon$  is a small positive number to prevent division equals to zero. The learning rate divided by the mean of gradients which is exponentially decreasing.

## A.3 Adam

Adam (Adaptive Moment Estimation) is another popular adaptive gradient descent algorithm, it proposed by [Kingma et al. \(2015\)](#). it combined both RMSprop and momentum, i.e. it not only compute and save the exponentially decaying of previous mean squared gradient, but also compute the previous momentum. More specifically,

$$\begin{aligned}m_t &= \beta_1 m_{t-1} + (1 - \beta_1) \mathbf{g}_t \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) \mathbf{g}_t^2\end{aligned}$$

---

<sup>5</sup>[Neural Networks for Machine Learning](#)

When  $m_t$  and  $v_t$  are initialized to zero vectors, [Kingma et al.](#) found that the they are unbiased. The final bias corrected terms are

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$
$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

combined both

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{\sqrt{\hat{v}} + \epsilon} \hat{m}_t$$

Where the default settings are  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 10^{-8}$ . In practice, Adam performs better than other adaptive gradient descent algorithm.

## A.4 leaky Relu

The leaky Relu is an non-linear activation, which is developed by [Maas et al. \(2013\)](#). Since the normal Relu is  $f(x) = \max(0, x)$ , when  $f(x)$  is negative, it's gradient will be zero. Therefore, if the output of  $f(x)$  mostly are zeros, the normal Relu can cause vanishing gradients. The leaky Relu is to solve this problem by multiply a positive number which is greater than 1, more specifically,

$$f(x) = \begin{cases} x & x \geq 0 \\ \frac{x}{a} & 0 < x \end{cases}$$

where  $a \in (1, +\infty)$ , thus, the leaky Relu slightly change the gradient of the negative part of  $f(x)$ 's output, and make the backpropagation not suffered by vanishing gradients.

# Bibliography

- [1] Sebastian Ruder, Ivan Vulić, and Anders Søgaard. [A Survey Of Cross-lingual Word Embedding Models](#). 2017.
- [2] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. [Efficient Estimation of Word Representations in Vector Space](#). 2013.
- [3] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. [Generative Adversarial Nets](#). 2014.
- [4] Yanghua Jin, Jiakai Zhang, Minjun Li, Yingtao Tian, Huachun Zhu, and Zhihao Fang. [Towards the Automatic Anime Characters Creation with Generative Adversarial Networks](#). 2017.
- [5] Christian Ledig, Lucas Theis, Ferenc Huszar, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and Wenzhe Shi. [Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network](#). 2016.
- [6] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. [Generating Videos with Scene Dynamics](#). 2016.
- [7] William Fedus, Ian Goodfellow, and Andrew M. Dai. [MaskGAN: Better Text Generation via Filling in the \\_\\_\\_\\_\\_](#). 2018.
- [8] Olof Mogren. [C-RNN-GAN: Continuous recurrent neural networks with adversarial training](#). 2016.
- [9] Scott Reed, Zeynep Akata, Xincheng Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. [Generative Adversarial Text to Image Synthesis](#). 2016.
- [10] Alexis Conneau, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. [Word Translation Without Parallel Data](#). 2018.
- [11] [Urysohn's Lemma - Wikipedia](#)
- [12] Chao Xing, Dong Wang, Chao Liu, and Yiye Lin. [Normalized Word Embedding and Orthogonal Transform for Bilingual Word Translation](#). 2015.
- [13] Martin Arjovsky, Soumith Chintala, and Léon Bottou. [Wasserstein GAN](#). 2017.

- [14] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. [Improved Training of Wasserstein GANs](#). 2017.
- [15] Xiang Wei, Boqing Gong, Zixia Liu, Wei Lu, Liqiang Wang. [Improving the Improved Training of Wasserstein GANs: A Consistency Term and Its Dual Effect](#). 2017.
- [16] Tomas Mikolov, Quoc V. Le, and Ilya Sutskever. [Exploiting Similarities among Languages for Machine Translation](#). 2017.
- [17] Cédric Villani. [Optimal transport, old and new](#). 2008.
- [18] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin and Aaron Courville. [Improved Training of Wasserstein GANs](#). 2017.
- [19] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. [Enriching Word Vectors with Subword Information](#). 2016.
- [20] Martin Arjovsky and Léon Bottou. [Towards Principled Methods for Training Generative Adversarial Networks](#). 2017.
- [21] Hariharan Narayanan and Sanjoy Mitter. [Sample complexity of testing the manifold hypothesis](#). 2010.
- [22] Miloš Radovanović, Alexandros Nanopoulos, and Mirjana Ivanović. [Hubs in Space: Popular Nearest Neighbors in High-Dimensional Data](#). 2010.
- [23] Jeff Johnson, Matthijs Douze, and Hervé Jégou. [Billion-scale similarity search with GPUs](#). 2017.
- [24] Sergey Ioffe and Christian Szegedy. [Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift](#). 2015.
- [25] Anders Søgaard, Sebastian Ruder and Ivan Vulić. [On the Limitations of Unsupervised Bilingual Dictionary Induction](#). 2018.
- [26] Jimmy Lei Ba, Jamie Ryan Kiros and Geoffrey E. Hinton. [Layer Normalization](#). 2016.
- [27] Sebastian Ruder. [An overview of gradient descent optimization algorithms](#). 2017.
- [28] Diederik P. Kingma and Jimmy Ba. [Adam: A Method for Stochastic Optimization](#). 2015.
- [29] Andrew L. Maas , Awni Y. Hannun and Andrew Y. Ng. [Rectifier Nonlinearities Improve Neural Network Acoustic Models](#). 2013
- [30] Laurens van der Maaten and Geoffrey Hinton [Visualizing data using t-SNE](#). 2008